IBM SPSS Modeler 18.2 - Guide des applications



Remarque

Avant d'utiliser le présent document et le produit associé, prenez connaissance des informations générales figurant à la section «Remarques», à la page 365.

Informations produit

LE PRESENT DOCUMENT EST LIVRE EN L'ETAT SANS AUCUNE GARANTIE EXPLICITE OU IMPLICITE. IBM DECLINE NOTAMMENT TOUTE RESPONSABILITE RELATIVE A CES INFORMATIONS EN CAS DE CONTREFACON AINSI QU'EN CAS DE DEFAUT D'APTITUDE A L'EXECUTION D'UN TRAVAIL DONNE.

Ce document est mis à jour périodiquement. Chaque nouvelle édition inclut les mises à jour. Les informations qui y sont fournies sont susceptibles d'être modifiées avant que les produits décrits ne deviennent eux-mêmes disponibles. En outre, il peut contenir des informations ou des références concernant certains produits, logiciels ou services non annoncés dans ce pays. Cela ne signifie cependant pas qu'ils y seront annoncés.

Pour plus de détails, pour toute demande d'ordre technique, ou pour obtenir des exemplaires de documents IBM, référez-vous aux documents d'annonce disponibles dans votre pays, ou adressez-vous à votre partenaire commercial.

Vous pouvez également consulter les serveurs Internet suivants :

- http://www.fr.ibm.com (serveur IBM en France)
- http://www.ibm.com/ca/fr (serveur IBM au Canada)
- http://www.ibm.com (serveur IBM aux Etats-Unis)

Compagnie IBM France Direction Qualité 17, avenue de l'Europe 92275 Bois-Colombes Cedex

© Copyright IBM France 2018. Tous droits réservés.

Cette édition s'applique à la version 18.2.0 d'IBM SPSS Modeler et à toutes les éditions et modifications ultérieures, sauf indication contraire dans les nouvelles éditions.

Table des matières

Avis aux lecteurs canadiens	vii
Chapitre 1. A propos d'IBM SPSS Modeler	. 1
Draduita IBM CDCC Madalar	' • 1
IDM CDCC Madalar	. 1
IDM SPSS Modeler	. 1
IDM SPSS Modeler Server	. 1
IBM SPSS Modeler Administration Console	. 2
IBM SPSS Modeler Batch	. 2
IBM SPSS Modeler Solution Publisher	. 2
Adaptateurs IBM SPSS Modeler Server pour IBM	~
SPSS Collaboration and Deployment Services .	. 2
Editions d'IBM SPSS Modeler.	. 2
Documentation	. 3
Documentation SPSS Modeler Professional	. 3
Documentation de SPSS Modeler Premium	. 4
Exemples d'application	. 4
Dossier Demos	. 4
Suivi des licences	. 5
Chapitre 2. Présentation du produit	. 7
Démarrage	. 7
Démarrage d'IBM SPSS Modeler	. 7
Lancement de l'application à partir de la ligne de	_
commande	. 7
Connexion à IBM SPSS Modeler Server	. 8
Connexion à Analytic Server	10
Changement de répertoire temporaire	11
Démarrage de plusieurs sessions IBM SPSS	
Modeler	. 11
Interface IBM SPSS Modeler en un clin d'oeil	12
Canevas de flux IBM SPSS Modeler	12
Palette de noeuds	13
Gestionnaires IBM SPSS Modeler	15
Projets IBM SPSS Modeler	16
Barre d'outils IBM SPSS Modeler	16
Personnalisation de la barre d'outils	18
Personnalisation de la fenêtre IBM SPSS Modeler	18
Modification de la taille des icônes d'un flux	19
Utilisation de la souris dans IBM SPSS Modeler	
Utilisation des touches de raccourci	20
	20 20
Impression	20 20 21
Impression .	20 20 21 22

Chapitre 3. Introduction à la

modélisation	,	•	•	•	-	•	. 23
Création du flux							. 24
Navigation dans le modèle							. 29
Evaluation du modèle							. 34
Scoring des enregistrements							. 37
Récapitulatif							. 37

Chapitre 4. Modélisation automatisée d'une cible indicateur	39
automatique)	20
Deprése d'historique	. 39
Création du flux	. 39
Creation du liux	. 40
Bécapitulatif	. 44
	. 49
Chapitre 5. Modélisation automatisée	
d'une cible continue	51
Valeurs de propriété (Numérisation automatique)	. 51
Données d'apprentissage	. 51
Création du flux	. 52
Comparaison des modèles	. 55
Récapitulatif	. 57
Chapitre 6. Préparation automatique de	
données (ADP)	59
Création du flux	59
Comparaison de l'exactitude des modèles	63
comparaison de rexactitude des modeles	. 05
Chapitre 7. Préparation des données	
pour l'analyse (Audit données)	67
Création du flux.	. 67
Navigation dans les statistiques et les graphiques	. 70
Traitement des valeurs éloignées et manquantes .	. 72
Chapitre 8. Traitements par	
médicaments (Granhiques	
avalerateiree/CE 0)	77
Lecture de données texte	. 77
Ajout d'une table	. 80
Création d'un graphique de distribution.	. 81
Création d'un nuage de points	. 82
Création d'un graphique Relations	. 83
Calcul d'un nouveau champ	. 85
Création d'un modèle	. 88
Navigation dans le modèle	. 90
Utilisation d'un noeud Analyse	91
Chapitre 9 Eiltrage des prédicteurs	. 71
Chapitre 3. I httage des predicteurs	. 71
(sélection de fonction)	. 91 93
(sélection de fonction)	93 . 93
(sélection de fonction)	93 . 93 . 96
(sélection de fonction)	93 . 93 . 96 . 97
(sélection de fonction)	93 . 93 . 96 . 97 . 99

Chapitre 10. Réduction de la longueur des chaînes de données d'entrée

(Noeud Recoder)							101
Réduction de la longueur des ch	naîr	nes	de	do	nne	ées	
d'entrée (Reclassifier)							. 101
Reclassification des données							. 101

Chapitre 11. Modélisation de la

réponse client (L	_is	te	de	e d	éc	isi	on	I)		107
Données d'historique										. 107
Création du flux .										. 108
Création du modèle										. 110
Calcul des mesures p	ers	onr	nali	sée	es a	veo	: E>	ccel	l	. 123
Modification du m	od	èle	Ex	cel						. 129
Enregistrement des ré	ésul	ltat	s							. 131

Chapitre 12. Classification des clients de télécommunications (régression

logistique multinomial	e).				133
Création du flux					. 133
Navigation dans le modèle.					. 136

logioliquo billolinulo,						•		
Création du flux								. 141
Navigation dans le modèle	•	•		•	•		•	. 147

. _ _

Chapitre 14. Prévision de l'utilisation de la bande passante (Séries

temporelles)	•			153
Prévision avec le noeud Séries temporel	les			. 153
Création du flux	•			. 154
Analyse des données	•			. 155
Définition des dates	•			. 158
Définition des cibles	•			. 160
Définition des intervalles de temps .				. 161
Création du modèle				. 162
Examen du modèle				. 164
Récapitulatif				. 171
Réapplication d'un modèle de séries ten	npoi	elle	2S	171
Récupération du flux	Ē.			. 172
Extraction du modèle enregistré				. 173
Génération d'un noeud de modélisat	ion			. 173
Génération d'un nouveau modèle .				. 173
Examen du nouveau modèle				. 174
Récapitulatif	•			. 177
=				

Chapitre 15. Prévision des ventes sur

catalogue (series	5 T	en	۱p	ore	elle	es)	•	•	•	•	179
Création du flux .											. 179
Analyse des données											. 182
Lissage exponentiel											. 182
ARIMA											. 187
Récapitulatif		•		•							. 191

Chapitre 16. Propositions aux clients Chapitre 17. Prévision des défauts de paiement (Réseau Bayésien) 205 Chapitre 18. Recyclage d'un modèle chaque mois (Réseau Bayésien) . . . 213 Chapitre 19. Campagne publicitaire (Réseau de neurones/Arbre C&RT) . . 223 Chapitre 20. Surveillance d'état (Réseau de neurones/C5.0) 227 Chapitre 21. Classification des clients de services de télécommunications Etude des résultats de l'utilisation de l'analyse discriminante pour classifier les clients de Chapitre 22. Analyse de données de survie avec censure par intervalle (modèles linéaires généralisés). . . . 245 Ajustement du modèle avec le traitement pour Réapparition prédite et probabilités de survie . . 253 Modélisation de la probabilité de réapparition par Tests des effets du modèle..</td Réapparition prédite et probabilités de survie . . . 265

Chapitre 23. Utilisation de la régression de Poisson pour analyser les taux de dommage aux navires

(modèles linéaires généralis	sés	5).	•	-	271
Ajustement d'une régression de Pois	sso	n			
"surdispersée"					. 271
Statistiques de qualité d'ajustement					. 276
Test composite					. 277
Tests des effets du modèle					. 277
Estimations des paramètres.					. 278
Ajustement des modèles alternatifs					. 279
Statistiques de qualité d'ajustement					. 280
Récapitulatif					. 281
Procédures apparentées					. 281
Lectures recommandées					. 281

Chapitre 24. Ajustement d'une

régression gamma à des déclarations de sinistre automobile (modèles

linéaires généralisés).		•			283
Création du flux						. 283
Estimations des paramètres	5.					. 287
Récapitulatif						. 287
Procédures apparentées .						. 288
Lectures recommandées .						. 288

Chapitre 25. Classification des

échantillons de cellu	les	s (S	SV	M)	-			289
Création du flux								. 290
Analyse des données								. 294
Essai d'une autre fonction								. 296
Comparaison des résultats								. 297
Récapitulatif		•				•	•	. 298

Observations censurees .	•	•	•	•	•	•	•	. 302
Codages de variables catég	gor	iell	es					. 303
Sélection des variables .								. 304
Moyennes des covariables								. 306

Courbe de	sur	vie										307
Courbe de	risq	ue										307
Evaluation	•	•										308
Suivi du nomł	ore	pré	vu	de	cli	ent	s re	eter	nus			312
Scoring												323
Récapitulatif.								•				327

Chapitre 27. Analyse d'un panier de

courses (Induction de règle/C5.0).	329
Accès aux données	329
Identification des analogies entre les articles du	
panier	331
Portrait des groupes d'acheteurs	334
Récapitulatif	335

Chapitre 28. Estimation des offres de

nouveaux véhicules (KNN).		-	337
Création du flux			. 338
Examen des sorties			. 342
Espace du prédicteur			. 343
Graphique des homologues			. 344
Tableau des voisins et des distances			. 346
Récapitulatif			. 346

Chapitre 29. Découverte des relations de causalité dans les métriques

métier (TCM)
Création du flux
Exécution d'une analyse
Graphique de qualité du modèle de système global 350
Système de modèle global
Diagrammes d'impact
Détermination des causes premières des valeurs
extrêmes
Exécution de scénario
Remarques
Marques
Dispositions relatives à la documentation du
produit
Index

Avis aux lecteurs canadiens

Le présent document a été traduit en France. Voici les principales différences et particularités dont vous devez tenir compte.

Illustrations

Les illustrations sont fournies à titre d'exemple. Certaines peuvent contenir des données propres à la France.

Terminologie

La terminologie des titres IBM peut différer d'un pays à l'autre. Reportez-vous au tableau ci-dessous, au besoin.

IBM France	IBM Canada
ingénieur commercial	représentant
agence commerciale	succursale
ingénieur technico-commercial	informaticien
inspecteur	technicien du matériel

Claviers

Les lettres sont disposées différemment : le clavier français est de type AZERTY, et le clavier français-canadien de type QWERTY.

OS/2 et Windows - Paramètres canadiens

Au Canada, on utilise :

- les pages de codes 850 (multilingue) et 863 (français-canadien),
- le code pays 002,
- le code clavier CF.

Nomenclature

Les touches présentées dans le tableau d'équivalence suivant sont libellées différemment selon qu'il s'agit du clavier de la France, du clavier du Canada ou du clavier des États-Unis. Reportez-vous à ce tableau pour faire correspondre les touches françaises figurant dans le présent document aux touches de votre clavier.

France	Canada	Etats-Unis
K (Pos1)	ĸ	Home
Fin	Fin	End
1 (PgAr)	\$	PgUp
(PgAv)	₹	PgDn
Inser	Inser	Ins
Suppr	Suppr	Del
Echap	Echap	Esc
Attn	Intrp	Break
Impr écran	ImpEc	PrtSc
Verr num	Num	Num Lock
Arrêt défil	Défil	Scroll Lock
(Verr maj)	FixMaj	Caps Lock
AltGr	AltCar	Alt (à droite)

Brevets

Il est possible qu'IBM détienne des brevets ou qu'elle ait déposé des demandes de brevets portant sur certains sujets abordés dans ce document. Le fait qu'IBM vous fournisse le présent document ne signifie pas qu'elle vous accorde un permis d'utilisation de ces brevets. Vous pouvez envoyer, par écrit, vos demandes de renseignements relatives aux permis d'utilisation au directeur général des relations commerciales d'IBM, 3600 Steeles Avenue East, Markham, Ontario, L3R 9Z7.

Assistance téléphonique

Si vous avez besoin d'assistance ou si vous voulez commander du matériel, des logiciels et des publications IBM, contactez IBM direct au 1 800 465-1234.

Chapitre 1. A propos d'IBM SPSS Modeler

IBM[®] SPSS Modeler est un ensemble d'outils d'exploration de données qui vous permet de développer rapidement, grâce à votre expertise métier, des modèles prédictifs et de les déployer dans des opérations métier afin de faciliter la prise de décision. Conçu autour d'un modèle confirmé, le modèle CRISP-DM, IBM SPSS Modeler prend en charge l'intégralité du processus d'exploration de données, des données à l'obtention de meilleurs résultats commerciaux.

IBM SPSS Modeler propose différentes méthodes de modélisation issues des domaines de l'apprentissage automatique, de l'intelligence artificielle et des statistiques. Les méthodes disponibles dans la palette Modélisation vous permettent d'extraire de nouvelles informations de vos données et de développer des modèles prédictifs. Chaque méthode possède ses propres avantages et est donc plus adaptée à certains types de problème spécifiques.

Il est possible d'acquérir SPSS Modeler comme produit autonome ou de l'utiliser en tant que client en combinaison avec SPSS Modeler Server. Plusieurs autres options sont également disponibles, telles que décrites dans les sections suivantes. Pour plus d'informations, voir https://www.ibm.com/analytics/us/en/technology/spss/.

Produits IBM SPSS Modeler

La famille des produits IBM SPSS Modeler et les logiciels associés sont composés des éléments suivants.

- IBM SPSS Modeler
- IBM SPSS Modeler Server
- IBM SPSS Modeler Administration Console (inclus avec IBM SPSS Deployment Manager)
- IBM SPSS Modeler Batch
- IBM SPSS Modeler Solution Publisher
- Adaptateurs IBM SPSS Modeler Server pour IBM SPSS Collaboration and Deployment Services

IBM SPSS Modeler

SPSS Modeler est une version complète du produit que vous installez et exécutez sur votre ordinateur personnel. Pour obtenir de meilleures performances lors du traitement de jeux de données volumineux, vous pouvez exécuter SPSS Modeler en mode local, comme produit autonome, ou l'utiliser en mode réparti, en association avec IBM SPSS Modeler Server.

Avec SPSS Modeler, vous pouvez créer des modèles prédictifs précis rapidement et de manière intuitive, sans aucune programmation. L'interface visuelle unique vous permet de visualiser facilement le processus d'exploration de données. Grâce aux analyses avancées intégrées au produit, vous pouvez découvrir des motifs et tendances masqués dans vos données. Vous pouvez modéliser les résultats et comprendre les facteurs qui les influencent, afin d'exploiter les opportunités commerciales et de réduire les risques.

SPSS Modeler est disponible en deux éditions : SPSS Modeler Professional et SPSS Modeler Premium. Pour plus d'informations, voir la rubrique «Editions d'IBM SPSS Modeler», à la page 2.

IBM SPSS Modeler Server

Grâce à une architecture client/serveur, SPSS Modeler adresse les demandes d'opérations très consommatrices de ressources à un logiciel serveur puissant. Il offre ainsi des performances accrues sur des jeux de données plus volumineux.

SPSS Modeler Server est un produit avec licence distincte qui s'exécute en permanence en mode d'analyse réparti sur un hôte de serveur en combinaison avec une ou plusieurs installations d'IBM SPSS Modeler. Ainsi, SPSS Modeler Server fournit des performances supérieures sur de grands jeux de données car les opérations nécessitant beaucoup de mémoire peuvent être effectuées sur le serveur sans télécharger de données sur l'ordinateur client. IBM SPSS Modeler Server prend également en charge l'optimisation SQL et propose des capacités de modélisation dans la base de données pour des performances et une automatisation améliorées.

IBM SPSS Modeler Administration Console

Modeler Administration Console est une interface graphique permettant de gérer de nombreuses options de SPSS Modeler Server qui peuvent également être configurées au moyen d'un fichier d'options. La console est incluse dans IBM SPSS Deployment Manager et peut être utilisée pour surveiller et configurer vos installations SPSS Modeler Server ; elle est disponible gratuitement pour les clients actuels de SPSS Modeler Server. L'application ne peut être installée que sur des ordinateurs Windows ; en revanche, elle peut administrer un serveur installé sur n'importe quelle plate-forme prise en charge.

IBM SPSS Modeler Batch

Alors que l'exploration de données est généralement un processus interactif, il est également possible d'exécuter SPSS Modeler à partir d'une ligne de commande sans recourir à l'interface utilisateur graphique. Par exemple, vous pouvez avoir des tâches longue durée ou répétitives à exécuter sans intervention de l'utilisateur. SPSS Modeler Batch est une version spécifique du produit qui prend en charge toutes les capacités d'analyse de SPSS Modeler sans avoir besoin d'accéder à l'interface utilisateur standard. SPSS Modeler Server est requis pour utiliser SPSS Modeler Batch.

IBM SPSS Modeler Solution Publisher

SPSS Modeler Solution Publisher est un outil qui permet de créer une version « packagée » d'un flux SPSS Modeler qui peut être exécutée par un moteur Runtime externe ou intégrée dans une application externe. Ainsi, vous pouvez publier et déployer des flux SPSS Modeler complets dans des environnements où SPSS Modeler n'est pas installé. SPSS Modeler Solution Publisher est fourni avec le service IBM SPSS Collaboration and Deployment Services - Scoring et nécessite une licence distincte. Avec cette licence, vous recevez SPSS Modeler Solution Publisher Runtime qui vous permet d'exécuter les flux publiés.

Pour plus d'informations sur SPSS Modeler Solution Publisher, reportez-vous à la documentation de IBM SPSS Collaboration and Deployment Services. Le Knowledge Center IBM SPSS Collaboration and Deployment Services contient des sections intitulées "IBM SPSS Modeler Solution Publisher" et "IBM SPSS Analytics Toolkit."

Adaptateurs IBM SPSS Modeler Server pour IBM SPSS Collaboration and Deployment Services

Différents adaptateurs pour IBM SPSS Collaboration and Deployment Services sont disponibles et permettent à SPSS Modeler et SPSS Modeler Server d'interagir avec un référentiel IBM SPSS Collaboration and Deployment Services. Ainsi, un flux SPSS Modeler déployé sur le référentiel peut être partagé par différents utilisateurs ou peut être accessible depuis l'application client léger IBM SPSS Modeler Advantage. Installez l'adaptateur sur le système qui héberge le référentiel.

Editions d'IBM SPSS Modeler

SPSS Modeler est disponible dans les éditions suivantes.

SPSS Modeler Professional

SPSS Modeler Professional offre tous les outils nécessaires à l'utilisation de la plupart des types de données structurées, tels que les comportements et interactions suivis dans les systèmes CRM, les

caractéristiques sociodémographiques, les comportements d'achat et les données de vente.

SPSS Modeler Premium

SPSS Modeler Premium est un produit avec licence distincte qui étend le champ d'applications de SPSS Modeler Professional afin de pouvoir traiter des données spécialisées et des données de texte non structurées. SPSS Modeler Premium inclut IBM SPSS Modeler Text Analytics :

IBM SPSS Modeler Text Analytics utilise des technologies linguistiques avancées et le traitement du langage naturel pour traiter rapidement une large variété de données textuelles non structurées, en extraire les concepts clés et les organiser pour les regrouper dans des catégories. Les concepts extraits et les catégories peuvent ensuite être combinés aux données structurées existantes, telles que les données démographiques, et appliqués à la modélisation grâce à la gamme complète d'outils d'exploration de données d'IBM SPSS Modeler, afin de favoriser une prise de décision précise et efficace.

IBM SPSS Modeler Subscription

IBM SPSS Modeler Subscription propose les mêmes fonctionnalités d'analyse prédictive que le client IBM SPSS Modeler traditionnel. L'édition Subscription vous permet de télécharger régulièrement des mises à jour de produit.

Documentation

Une documentation est disponible dans le menu Aide de SPSS Modeler. Elle ouvre le Knowledge Center, qui est disponible au public en dehors du produit.

La documentation complète de chaque produit (y compris les instructions d'installation) au format PDF est également disponible dans un dossier compressé distinct, avec le téléchargement du produit. Les documents PDF peuvent également être téléchargés depuis le Web sur le site http://www.ibm.com/support/docview.wss?uid=swg27046871.

Documentation SPSS Modeler Professional

La suite de documentation SPSS Modeler Professional (à l'exception des instructions d'installation) est la suivante.

- IBM SPSS Modeler Guide d'utilisation. Introduction générale à SPSS Modeler : création de flux de données, traitement des valeurs manquantes, création d'expressions CLEM, utilisation de projets et de rapports, et regroupement des flux pour le déploiement dans IBM SPSS Collaboration and Deployment Services ou IBM SPSS Modeler Advantage.
- Noeuds source, de processus et de sortie IBM SPSS Modeler. Descriptions de tous les noeuds utilisés pour lire, traiter et renvoyer les données de sortie dans différents formats. En pratique, cela signifie tous les noeuds autres que les noeuds de modélisation.
- Noeuds Modélisation IBM SPSS Modeler. Descriptions de tous les noeuds utilisés pour créer des modèles d'exploration de données. IBM SPSS Modeler propose différentes méthodes de modélisation issues des domaines de l'apprentissage automatique, de l'intelligence artificielle et des statistiques.
- **IBM SPSS Modeler Guide des applications.** Les exemples de ce guide fournissent des introductions brèves et ciblées aux méthodes et techniques de modélisation. Une version en ligne de ce guide est également disponible dans le menu Aide. Pour plus d'informations, voir la rubrique «Exemples d'application», à la page 4.
- **IBM SPSS Modeler Guide de génération de scripts Python et d'automatisation.** Ce manuel fournit des informations sur l'automatisation du système via des scripts Python, notamment sur les propriétés pouvant être utilisées pour manipuler les noeuds et les flux.
- **IBM SPSS Modeler Guide de déploiement.** Informations sur l'exécution des flux IBM SPSS Modeler comme étapes des travaux d'exécution sous IBM SPSS Deployment Manager.

- IBM SPSS Modeler CLEF Guide du développeur. CLEF permet d'intégrer des programmes tiers tels que des programmes de traitement de données ou des algorithmes de modélisation en tant que noeuds dans IBM SPSS Modeler.
- **IBM SPSS Modeler Guide d'exploration de base de données.** Informations sur la manière de tirer parti de la puissance de votre base de données pour améliorer les performances et étendre la gamme des capacités d'analyse via des algorithmes tiers.
- **IBM SPSS Modeler Server Guide d'administration et des performances.** Informations sur le mode de configuration et d'administration d'IBM SPSS Modeler Server.
- **IBM SPSS Deployment Manager Guide d'utilisation.** Informations sur l'utilisation de l'interface utilisateur de la console d'administration incluses dans l'application Deployment Manager pour la surveillance et la configuration d'IBM SPSS Modeler Server.
- **IBM SPSS Modeler Guide CRISP-DM.** Guide détaillé sur l'utilisation de la méthodologie CRISP-DM pour l'exploration de données avec SPSS Modeler
- **IBM SPSS Modeler Batch Guide d'utilisation.** Guide complet sur l'utilisation d'IBM SPSS Modeler en mode de traitement par lots, avec des détails sur l'exécution en mode de traitement par lots et les arguments de ligne de commande. Ce guide est disponible au format PDF uniquement.

Documentation de SPSS Modeler Premium

La suite de documentation SPSS Modeler Premium (à l'exception des instructions d'installation) est la suivante.

• Guide d'utilisation de SPSS Modeler Text Analytics . Informations sur l'utilisation des analyses de texte avec SPSS Modeler, notamment sur les noeuds Text Mining, l'espace de travail interactif, les modèles et d'autres ressources.

Exemples d'application

Tandis que les outils d'exploration de données de SPSS Modeler peuvent vous aider à résoudre une grande variété de problèmes métier et organisationnels, les exemples d'application fournissent des introductions brèves et ciblées aux méthodes et aux techniques de modélisation. Les jeux de données utilisés ici sont beaucoup plus petits que les énormes entrepôts de données gérés par certains Data miners, mais les concepts et les méthodes impliqués peuvent être adaptés à des applications réelles.

Pour accéder aux exemples, cliquez sur Exemples d'application dans le menu Aide de SPSS Modeler.

Les fichiers de données et les flux d'échantillons sont installés dans le dossier Demos, sous le répertoire d'installation du produit. Pour plus d'informations, voir «Dossier Demos».

Exemples de modélisation de bases de données. Consultez les exemples dans le document *IBM SPSS Modeler Guide d'exploration de base de données.*

Exemples de génération de scripts. Consultez les exemples dans le document *IBM SPSS Modeler Guide de génération de scripts et d'automatisation*.

Dossier Demos

Les fichiers de données et les flux d'échantillons utilisés avec les exemples d'application sont installés dans le dossier Demos, sous le répertoire d'installation du produit (par exemple : C:\Program Files\IBM\SPSS\Modeler\<version>\Demos). Ce dossier est également accessible à partir du groupe de programmes IBM SPSS Modeler, dans le menu Démarrer de Windows ou en cliquant sur Demos dans la liste des répertoires récents de la boîte de dialogue **Fichier > Ouvrir un flux**.

Suivi des licences

Lorsque vous utilisez SPSS Modeler, l'utilisation des licences est suivie et consignée à intervalles réguliers. Les métriques de licence consignées sont *AUTHORIZED_USER* et *CONCURRENT_USER* et le type de métrique consigné dépend du type de licence dont vous disposez pour SPSS Modeler.

Les fichiers journaux générés peuvent être traités par IBM License Metric Tool, à partir duquel vous pouvez générer des rapports d'utilisation de licence.

Les fichiers journaux des licences sont créés dans le répertoire dans lequel les fichiers journaux de SPSS Modeler Client sont enregistrés (par défaut, %ALLUSERSPROFILE%/IBM/SPSS/Modeler/<version>/log).

Chapitre 2. Présentation du produit

Démarrage

En tant qu'application d'exploration de données, IBM SPSS Modeler constitue une méthode stratégique de recherche de relations utiles dans les grands jeux de données. Contrairement aux méthodes statistiques plus traditionnelles, il n'est pas indispensable de savoir ce que vous recherchez avant de commencer. Vous pouvez explorer vos données, créer divers modèles et explorer diverses relations, jusqu'à ce que vous trouviez des informations utiles.

Démarrage d'IBM SPSS Modeler

Pour démarrer l'application, cliquez :

Démarrer > [Tous] Programmes > IBM SPSS Modeler <version> > IBM SPSS Modeler <version>

La fenêtre principale apparaît après quelques secondes.



Figure 1. Fenêtre d'application principale IBM SPSS Modeler

Lancement de l'application à partir de la ligne de commande

Vous pouvez utiliser la ligne de commande de votre système d'exploitation pour lancer IBM SPSS Modeler comme suit :

- 1. Dans le cas d'un ordinateur sur lequel est installé IBM SPSS Modeler, ouvrez une fenêtre DOS ou une invite de commande.
- 2. Pour lancer l'interface IBM SPSS Modeler en mode interactif, tapez la commande modelerclient suivie des arguments requis, par exemple :

modelerclient -stream report.str -execute

Les arguments disponibles (indicateurs) vous permettent de vous connecter à un serveur, de charger des flux, d'exécuter des scripts ou d'indiquer les autres paramètres nécessaires.

Connexion à IBM SPSS Modeler Server

Il est possible d'exécuter IBM SPSS Modeler comme une application autonome ou un comme un client connecté directement à IBM SPSS Modeler Server ou à IBM SPSS Modeler Server or à un cluster de serveurs par le biais du Coordinateur des processus connecté à partir de IBM SPSS Collaboration and Deployment Services. Le statut de la connexion apparaît en bas à gauche de la fenêtre IBM SPSS Modeler.

Lorsque vous souhaitez vous connecter à un serveur, vous pouvez saisir manuellement son nom ou sélectionner un nom que vous aurez préalablement défini. En revanche, si vous avez IBM SPSS Collaboration and Deployment Services, vous avez la possibilité de chercher dans une liste de serveurs ou de clusters de serveurs à partir de la boîte de dialogue Connexion au serveur. Vous pouvez naviguer via les services Statistiques s'exécutant sur un réseau grâce au Coordinateur des processus.

Pour vous connecter à un serveur

- Dans le menu Outils, cliquez sur Connexion au serveur. La boîte de dialogue Connexion au serveur s'affiche. Vous pouvez également cliquer deux fois sur la zone d'état de la connexion dans la fenêtre IBM SPSS Modeler.
- 2. Dans la boîte de dialogue, indiquez les options de connexion à l'ordinateur du serveur local ou sélectionnez une connexion dans le tableau.
 - Cliquez sur **Ajouter** ou **Modifier** pour ajouter ou modifier une connexion. Pour plus d'informations, voir la rubrique «Ajout et modification d'une connexion à IBM SPSS Modeler Server», à la page 9.
 - Cliquez sur **Rechercher** pour accéder au serveur ou à un cluster de serveurs dans le Coordinateur de processus. Pour plus d'informations, voir la rubrique «Recherche de serveurs dans IBM SPSS Collaboration and Deployment Services», à la page 9.

Tableau de serveur. Ce tableau comprend un ensemble de connexions au serveur définies. Il affiche la connexion par défaut, le nom du serveur, sa description et le numéro du port. Vous pouvez ajouter manuellement une nouvelle connexion ainsi que sélectionner ou rechercher une connexion existante. Pour définir un serveur particulier comme connexion par défaut, cochez la case dans la colonne Par défaut du tableau de la connexion.

Chemin de données par défaut. Indiquez le chemin d'accès aux données situées sur l'ordinateur serveur. Cliquez sur le bouton ... pour accéder à l'emplacement requis.

Définir les données d'identification. Laissez cette case décochée pour permettre à la fonction de **connexion unique** de se connecter au serveur à l'aide de vos informations de nom d'utilisateur et de mot de passe locaux. Si la connexion unique n'est pas disponible, ou si vous cochez la case pour désactiver la connexion unique (par exemple pour vous connecter à un compte administrateur), les champs suivants sont activés et vous permettent d'entrer vos informations d'identification.

ID utilisateur. Entrez le nom d'utilisateur avec lequel effectuer la connexion au serveur.

Mot de passe. Entrez le mot de passe associé au nom d'utilisateur défini.

Domaine. Indiquez le domaine utilisé pour la connexion au serveur. Le nom de domaine n'est requis que si l'ordinateur serveur se trouve dans un autre domaine Windows que l'ordinateur client.

3. Cliquez sur OK pour terminer la connexion.

Pour se déconnecter d'un serveur

- Dans le menu Outils, cliquez sur Connexion au serveur. La boîte de dialogue Connexion au serveur s'affiche. Vous pouvez également cliquer deux fois sur la zone d'état de la connexion dans la fenêtre IBM SPSS Modeler.
- 2. Dans la boîte de dialogue, sélectionnez le serveur local, puis cliquez sur OK.

Ajout et modification d'une connexion à IBM SPSS Modeler Server

Dans la boîte de dialogue Connexion au serveur, vous pouvez modifier ou ajouter une connexion au serveur. Cliquez sur Ajouter pour accéder à une boîte de dialogue Ajouter/Modifier un serveur non renseignée, dans laquelle vous pourrez entrer les données de la connexion au serveur. Si vous sélectionnez une connexion existante et cliquez sur Modifier, dans la boîte de dialogue Connexion au serveur, la boîte de dialogue Ajouter/Modifier un serveur s'ouvre, affichant les données de cette connexion, vous permettant ainsi d'apporter toutes les modifications que vous souhaitez.

Remarque : Vous ne pouvez pas modifier une connexion au serveur qui a été ajoutée à partir de IBM SPSS Collaboration and Deployment Services, car le nom, le port et d'autres détails sont définis dans IBM SPSS Collaboration and Deployment Services. Les pratiques recommandées indiquent d'utiliser les mêmes ports pour communiquer avec IBM SPSS Collaboration and Deployment Services et avec SPSS Modeler Client. Ces ports peuvent être définis via les paramètres max_server_port et min_server_port dans le fichier options.cfg.

Pour ajouter des connexions au serveur

- 1. Dans le menu Outils, cliquez sur **Connexion au serveur**. La boîte de dialogue Connexion au serveur s'affiche.
- 2. Dans la boîte de dialogue, cliquez sur **Ajouter**. La boîte de dialogue Ajouter/Modifier un serveur pour la connexion au serveur s'ouvre.
- **3**. Saisissez les données de connexion au serveur puis cliquez sur **OK** pour enregistrer la connexion et retourner à la boîte de dialogue Connexion au serveur.
- **Serveur.** Indiquez un serveur disponible ou sélectionnez-en un dans la liste. L'ordinateur serveur peut être identifié par un nom alphanumérique (par exemple, *monserveur*) ou une adresse IP qui lui est affectée (par exemple, 202.123.456.78).
- **Port.** Indiquez le numéro de port d'écoute du serveur. Si ce port par défaut ne fonctionne pas, demandez à l'administrateur système le numéro de port correct.
- Description. Saisissez une description optionnelle pour la connexion à ce serveur.
- Coder la connexion (utiliser SSL). Indique si une connexion SSL (Secure Sockets Layer) doit être utilisée. Le protocole SSL est fréquemment utilisé pour la sécurisation des données sur un réseau. Pour pouvoir utiliser cette fonction, vous devez activer le protocole SSL sur le serveur hébergeant IBM SPSS Modeler Server. Si nécessaire, contactez votre administrateur local pour plus d'informations.

Pour modifier des connexions au serveur

- 1. Dans le menu Outils, cliquez sur **Connexion au serveur**. La boîte de dialogue Connexion au serveur s'affiche.
- 2. Dans la boîte de dialogue, sélectionnez la connexion que vous souhaitez modifier puis cliquez sur **Modifier**. La boîte de dialogue Ajouter/Modifier un serveur pour la connexion au serveur s'ouvre.
- **3**. Modifiez les données de connexion au serveur puis cliquez sur **OK** pour enregistrer les changements et retourner à la boîte de dialogue Connexion au serveur.

Recherche de serveurs dans IBM SPSS Collaboration and Deployment Services

Au lieu d'entrer manuellement une connexion au serveur, vous pouvez sélectionner un serveur ou un cluster de serveurs disponible sur le réseau par le biais du Coordinateur de processus, disponible dans IBM SPSS Collaboration and Deployment Services. Un cluster de serveurs contient plusieurs serveurs, et permet au Coordinateur de processus de déterminer le serveur qui répond le mieux à la demande de traitement.

Bien que vous ne vous puissiez pas ajouter manuellement de serveurs dans la boîte de dialogue Connexion au serveur, la recherche de serveurs disponibles vous permet de vous connecter aux serveurs sans que vous ayez besoin de connaître le nom du serveur et le numéro du port. Ces informations sont fournies automatiquement. Il vous faut néanmoins corriger les données de connexion telles que le nom de l'utilisateur, le domaine et le mot de passe.

Remarque : Si vous n'avez pas accès au Coordinateur de processus, vous pouvez tout de même saisir manuellement le nom du serveur auquel vous souhaitez vous connecter ou sélectionner un nom que vous aurez défini au préalable. Pour plus d'informations, voir la rubrique «Ajout et modification d'une connexion à IBM SPSS Modeler Server», à la page 9.

Pour rechercher des serveurs et des clusters de serveurs

- 1. Dans le menu Outils, cliquez sur **Connexion au serveur**. La boîte de dialogue Connexion au serveur s'affiche.
- 2. Cliquez sur **Rechercher** pour ouvrir la boîte de dialogue Recherche de serveurs. Si vous n'êtes plus connecté à IBM SPSS Collaboration and Deployment Services, lors de votre tentative d'accès au Coordinateur de processus, il vous sera demandé de vous reconnecter.
- 3. Sélectionnez le serveur ou le cluster de serveurs dans la liste.
- 4. Cliquez sur **OK** pour fermer la boîte de dialogue et ajouter cette connexion au tableau de la boîte de dialogue Connexion au serveur.

Connexion à Analytic Server

Si plusieurs serveurs Analytic Server sont disponibles, vous pouvez utiliser la boîte de dialogue Connexion à Analytic Server pour définir plusieurs serveurs à utiliser dans IBM SPSS Modeler. Votre administrateur a peut-être déjà configuré un serveur Analytic Server par défaut dans le fichier <chemin_installation_Modeler>/config/options.cfg. Toutefois, vous pouvez utiliser d'autres serveurs disponibles après les avoir définis. Par exemple, si vous utilisez les noeuds Source et Exportation d'Analytic Server, vous pouvez utiliser des connexions Analytic Server dans différentes branches d'un flux afin que lorsque chaque branche est exécutée, elle utilise son propre serveur Analytic Server, sans extraction de données dans IBM SPSS Modeler Server. Notez que si une branche contient plusieurs connexions Analytic Server, les données seront extraites des serveurs Analytic Server vers IBM SPSS Modeler Server. Pour plus d'informations, notamment sur les restrictions, voir Propriétés du flux Analytic Server.

Pour créer une connexion Analytic Server, accédez à **Outils** > **Connexions Analytic Server** et spécifiez les informations requises dans les sections ci-après de la boîte de dialogue.

Connexion

URL. Entrez l'URL d'Analytic Server au format https://nomhôte:port/racinecontexte, nomhôte correspondant à l'adresse IP ou au nom d'hôte d'Analytic Server, port, au numéro de port, et racinecontexte, à la racine de contexte d'Analytic Server.

Locataire. Entrez le nom du locataire dont IBM SPSS Modeler Server est membre. Si vous ne connaissez pas le locataire, contactez votre administrateur.

Authentification

Mode. Sélectionnez l'un des modes d'authentification ci-après.

- Nom d'utilisateur et mot de passe. Vous devez entrer le nom d'utilisateur et le mot de passe.
- Données d'identification stockées. Vous devez sélectionner des informations d'identification dans IBM SPSS Collaboration and Deployment Services Repository.
- Kerberos. Vous devez entrer le nom principal du service et le chemin d'accès au fichier de configuration. Si vous ne connaissez pas ces informations, contactez votre administrateur.

Nom d'utilisateur. Entrez le nom d'utilisateur d'Analytic Server.

Mot de passe. Entrez le mot de passe d'Analytic Server.

Connecter. Cliquez sur Connecter pour tester la nouvelle connexion.

Connexions

Une fois que vous avez spécifié les informations ci-dessus et cliqué sur **Connecter**, la connexion est ajoutée à cette table Connexions. Si vous devez supprimer une connexion, sélectionnez-la et cliquez sur **Supprimer**.

Si votre administrateur a défini une connexion Analytic Server par défaut dans le fichier options.cfg, vous pouvez cliquer sur **Ajouter une connexion par défaut** pour l'ajouter également à vos connexions disponibles. Vous serez invité à entrer le nom d'utilisateur et le mot de passe.

Changement de répertoire temporaire

Certaines opérations effectuées par IBM SPSS Modeler Server peuvent nécessiter la création de fichiers temporaires. Par défaut, IBM SPSS Modeler crée les fichiers temporaires dans le répertoire temporaire du système. Vous pouvez modifier l'emplacement du répertoire temporaire en effectuant les opérations suivantes.

- 1. Créez un répertoire intitulé spss et un sous-répertoire intitulé servertemp.
- Editez le fichier options.cfg, situé dans le répertoire /config du dossier d'installation d'IBM SPSS Modeler. Editez le paramètre temp_directory de ce fichier en saisissant : temp_directory, "C:/spss/servertemp".
- **3**. Redémarrez ensuite le service IBM SPSS Modeler Server. Pour ce faire, cliquez sur **Services** dans les outils d'administration du Panneau de configuration Windows. Il vous suffit d'arrêter le service et de le redémarrer pour appliquer les modifications apportées. Vous pouvez redémarrer l'ordinateur pour redémarrer le service.

Tous les fichiers temporaires sont désormais écrits dans ce nouveau répertoire.

Remarque :

- Des barres obliques doivent être utilisées.
- Le paramètre temp_directory n'est pas applicable lors de l'exécution de flux d'évaluation via des travaux IBM SPSS Collaboration and Deployment Services. Lorsque vous exécutez ce type de travail, un fichier temporaire est créé. Par défaut, le fichier est enregistré dans le répertoire d'installation d'IBM SPSS Modeler Server. Vous pouvez changer le dossier de données par défaut dans lequel les fichiers temporaires sont enregistrés lorsque vous créez la connexion IBM SPSS Modeler Server dans IBM SPSS Modeler.

Démarrage de plusieurs sessions IBM SPSS Modeler

Si vous devez lancer plus d'une session IBM SPSS Modeler à la fois, vous devez effectuer certaines modifications de vos paramètres IBM SPSS Modeler et Windows. Par exemple, il vous faudra effectuer ces modifications si vous avez deux licences de serveur distinctes et que vous souhaitez exécuter deux flux pour deux serveurs distincts du même ordinateur client.

Pour activer plusieurs sessions IBM SPSS Modeler :

1. Cliquez sur :

Démarrer > [Tous les] Programmes > IBM SPSS Modeler

- 2. Dans le raccourci d'IBM SPSS Modeler (celui avec l'icône), cliquez avec le bouton droit de la souris et sélectionnez **Propriétés**.
- 3. Dans la zone de texte Cible, ajoutez -noshare à la fin de la chaîne.

- Dans Windows Explorer, sélectionnez : Outils > Options des dossiers...
- 5. Dans l'onglet Types de fichier, sélectionnez l'option Flux IBM SPSS Modeler et cliquez sur Avancé.
- 6. Dans la boîte de dialogue Modifier le type de fichier, sélectionnez Ouvrir avec IBM SPSS Modeler et cliquez sur **Edition**.
- 7. Dans la zone de texte **Application utilisée pour effectuer l'action**, ajoutez -noshare avant l'argument -flux.

Interface IBM SPSS Modeler en un clin d'oeil

A chaque étape du processus d'exploration des données, l'interface simplifiée IBM SPSS Modeler sollicite votre expertise métier. Les algorithmes de modélisation, comme la prévision, la classification, la segmentation et la détection d'association, permettent l'obtention de modèles performants et précis. Les résultats du modèle peuvent facilement être déployés et lus dans des bases de données, dans IBM SPSS Statistics et dans de nombreuses autres applications.

Travailler avec IBM SPSS Modeler est un processus en trois étapes de travail avec les données.

- Pour commencer, lisez les données d'IBM SPSS Modeler.
- Ensuite, exécutez les données par une série de manipulations.
- Et pour finir, envoyez les données vers une destination choisie.

Cette séquence d'opérations est appelée **flux de données** car les données circulent, enregistrement par enregistrement, de la source à la destination (modèle ou type de sortie de données), en passant par chaque manipulation.



Figure 2. Un flux simple

Canevas de flux IBM SPSS Modeler

Le canevas de flux est la plus grande zone de la fenêtre IBM SPSS Modeler. C'est dans cette zone que vous créez et manipulez les flux de données.



Figure 3. Espace de travail d'IBM SPSS Modeler (vue par défaut)

Les flux sont créés en dessinant des diagrammes des opérations de données nécessaires à votre entreprise sur l'espace de travail principal de l'interface. Chaque opération est représentée par une icône ou un **noeud**, et les noeuds sont reliés entre eux dans un **flux** représentant le flux de données passant par chaque opération.

Vous pouvez utiliser plusieurs flux à la fois dans IBM SPSS Modeler, que ce soit dans le même canevas de flux ou par l'ouverture d'un nouveau flux. Au cours d'une session, les flux sont stockés dans le gestionnaire de flux, en haut à droite de la fenêtre IBM SPSS Modeler.

Remarque : Si vous utilisez un MacBook avec le paramètre de pavé tactile **Force Click and haptic feedback** activé, le fait de glisser-déplacer la palette de noeuds sur le canevas de flux peut entraîner l'ajout de noeuds en double au canevas. Pour éviter ce problème, nous vous recommandons de désactiver la préférence système de pavé tactile **Force Click and haptic feedback**.

Palette de noeuds

La plupart des données et des outils de modélisation de SPSS Modeler sont disponibles dans la *Palette de noeuds*, au bas de la fenêtre sous l'espace de travail de flux.

Par exemple, l'onglet de la palette **Ops. sur lignes** contient des noeuds permettant d'effectuer des opérations sur les *lignes*, telles que la sélection, la fusion et l'ajout.

Pour ajouter des noeuds à l'espace de travail, double-cliquez sur les icônes de la palette de noeuds ou faites glisser les icônes vers l'espace de travail. Vous pouvez ensuite les relier afin de créer un *flux*

représentant le flux des données.



Figure 4. Onglet Ops sur enregistrements de la palette de noeuds

Chaque onglet de palette contient un ensemble de noeuds connexes employés pour différentes étapes des opérations de flux, comme :

- Les noeuds **Sources** amènent les données dans SPSS Modeler.
- Les noeuds **Ops sur lignes** sont utilisés pour les opérations sur les *lignes* de données, comme la sélection, la fusion et l'ajout.
- Les noeuds **Ops sur champs** sont utilisés pour les opérations sur les *champs* de données, comme le filtrage, le calcul de nouveaux champs et la détermination du niveau de mesure de champs particuliers.
- Les noeuds **Graphiques** sont utilisés pour visualiser les données avant et après la modélisation. Les graphiques peuvent être des tracés, des histogrammes, des noeuds relations, ainsi que des graphiques d'évaluation.
- Les noeuds de **Modélisation** utilisent les algorithmes de modélisation disponibles dans SPSS Modeler, tel que les réseaux neuronaux, les arbres de décisions, les algorithmes de groupement, et l'organisation de données.
- Les noeuds de **Modélisation de la base de données** utilisent les algorithmes de modélisation disponibles dans les bases de données Microsoft SQL Server, IBM Db2, Oracle et Netezza.
- Les noeuds de **Sortie** produisent diverses sorties pour les données, les graphiques et les résultats de modèles qui peuvent être affichés dans SPSS Modeler.
- Les noeuds **Exporter** produisent diverses sorties qui peuvent être affichées dans des applications externes telles que IBM SPSS Data Collection ou Excel.
- Les noeuds **IBM SPSS Statistics** importent des données à partir de, ou exportent des données vers IBM SPSS Statistics, et exécutent aussi des procédures IBM SPSS Statistics.
- Les noeuds **Python** peuvent être utilisés pour exécuter des algorithmes Python.
- · Les noeuds Spark peuvent être utilisés pour exécuter des algorithmes Spark.

Au fur et à mesure que vous maîtrisez mieux l'application SPSS Modeler, vous pouvez personnaliser le contenu de la palette en fonction de vos besoins.

A gauche de la palette des noeuds, vous pouvez filtrer les noeuds qui s'affichent en sélectionnant Supervisé, Association ou Segmentation.

Situé sous la palette de noeuds, un panneau de rapports fournit des informations sur la progression des diverses opérations, telles que la lecture des données dans le flux de données. Egalement situé sous la palette de noeuds, un panneau de statut fournit des informations sur l'activité actuelle de l'application, ainsi que des indications lorsqu'une saisie par l'utilisateur est requise.

Remarque : Si vous utilisez un MacBook avec le paramètre de pavé tactile **Force Click and haptic feedback** activé, le fait de glisser-déplacer la palette de noeuds sur le canevas de flux peut entraîner l'ajout de noeuds en double au canevas. Pour éviter ce problème, nous vous recommandons de désactiver la préférence système de pavé tactile **Force Click and haptic feedback**.

Gestionnaires IBM SPSS Modeler

En haut à droite de la fenêtre se trouve le panneau des gestionnaires. Il contient trois onglets qui permettent de gérer les flux, les sorties et les modèles.

Vous pouvez utiliser l'onglet Flux pour ouvrir, renommer, enregistrer et supprimer les flux créés dans une session.

drug			
- druaplo	t		
druarer	oort		
fraud			
	drug drugplo drugrep <mark>fraud</mark>	drug drugplot drugreport <mark>fraud</mark>	drug drugplot drugreport fraud

Figure 5. Onglet Flux



Figure 6. Onglet Sorties

L'onglet Sorties contient différents fichiers, tels que des graphiques et des tableaux, produits par des opérations de flux dans IBM SPSS Modeler. Vous pouvez afficher, enregistrer, renommer et fermer les tableaux, les graphiques et les rapports qui figurent dans cet onglet.

N	1	>	
Drug	olain		
Didg	cram	Ivalue	
3			
Drug			

Figure 7. Onglet Modèles qui contient des nuggets de modèles

L'onglet Modèle est le plus puissant des onglets du gestionnaire. Cet onglet contient tous les **nuggets** de modèle qui contiennent les modèles générés dans IBM SPSS Modeler, pour la session en cours. Vous pouvez accéder à ces modèles directement à partir de l'onglet Modèles ou les ajouter au flux dans l'espace de travail.

Projets IBM SPSS Modeler

Dans la partie inférieure droite de la fenêtre se trouve le panneau de projet qui permet de créer et de gérer les **projets** d'exploration de données (groupes de fichiers en rapport avec une tâche d'exploration de données). Vous pouvez afficher les projets créés de deux façons dans IBM SPSS Modeler : dans la vue Classes et dans la vue CRISP-DM.



Figure 8. Vue CRISP-DM

L'onglet CRISP-DM permet d'organiser les projets en fonction de la méthodologie Cross-Industry Standard Process for Data Mining commune utilisée dans le domaine. Que vous soyez un utilisateur chevronné ou novice, l'outil CRISP vous aidera à mieux organiser et communiquer vos efforts.

CRISP-DM Classes
 (unsaved project) Streams Orders Cenerated Models Tables, Graphs & Reports Patient Records (8 fields, 200 records) Distribution of Drug Other

Figure 9. Vue Classes

L'onglet Classes permet d'organiser votre travail dans IBM SPSS Modeler en catégories, selon les types d'objet que vous créez. Cette vue est utile lorsque vous effectuez l'inventaire des données, des flux et des modèles.

Barre d'outils IBM SPSS Modeler

Une barre d'outils, composée d'icônes fournissant des options très utiles, se trouve en haut de la fenêtre IBM SPSS Modeler. Voici les boutons de la barre d'outils et leurs fonctions.



Créer un flux



Permet d'ouvrir un flux



Enregistrer le flux

Imprimer le flux actuel



Déplacer la sélection vers le Presse-papiers



Copier dans le Presse-papiers



Coller le contenu du Presse-papiers dans la sélection



Annuler la dernière action



Rétablir la dernière action



Recherche de noeuds



Editer les propriétés du flux



Aperçu de génération SQL

Exécuter la sélection de flux



Exécuter le flux actuel



Ajouter un super noeud

Insérer un commentaire



Zoom avant (super noeuds uniquement)

Arrêter le flux (actif uniquement

pendant l'exécution du flux)



Aucun balisage dans le flux



Masquer le balisage de flux (le cas échéant)



Afficher le balisage de flux masqué

Zoom arrière (super noeuds uniquement)



Ouvrir un flux dans IBM SPSS Modeler Advantage

Le balisage de flux se compose des commentaires de flux, des liens de modèle et des indications de branche de scoring.

Les liens de modèle sont décrits dans le guide Noeuds de modélisation IBM SPSS.

Personnalisation de la barre d'outils

Vous pouvez modifier plusieurs aspects de la barre d'outils, tels que :

- choisir si elle sera affichée ou non
- · Choisir si les icônes comporteront ou non des info-bulles
- · Choisir si elle utilisera des petites ou des grandes icônes

activer ou désactiver l'affichage de la barre d'outils :

```
    Dans le menu principal, cliquez sur :

        Vue > Barre d'outils > Afficher
```

Pour modifier les paramètres des info-bulles ou de la taille des icônes :

1. Dans le menu principal, cliquez sur :

```
Vue > Barre d'outils > Personnaliser
```

Cliquez sur Afficher les info-bulles ou Gros boutons le cas échéant.

Personnalisation de la fenêtre IBM SPSS Modeler

Vous pouvez utiliser les séparateurs situés entre les différentes zones de l'interface SPSS Modeler pour redimensionner ou fermer des outils en fonction de vos besoins. Par exemple, si vous travaillez avec un flux volumineux, vous pouvez utiliser les petites flèches situées sur chaque séparateur pour fermer la palette de noeuds, le panneau des gestionnaires et le panneau des projets. Ainsi, vous agrandissez le canevas de flux et libérez suffisamment d'espace pour les flux volumineux ou multiples.

A partir du menu Vue, vous pouvez aussi cliquer sur **Palette de noeuds**, **Gestionnaires** ou **Projet** pour activer ou désactiver l'affichage de ces éléments.



Figure 10. Canevas de flux agrandi

Vous pouvez également garder ouvertes la palette des noeuds, et les panneaux des gestionnaires et des projets, et utiliser les barres de défilement du canevas de flux pour vous déplacer dans cet espace ; ces barres sont situées sur le côté et en bas de la fenêtre SPSS Modeler.

Vous pouvez aussi commander l'affichage du balisage de l'écran, lequel se compose des commentaires de flux, des liens de modèles et des indications de branche de scoring. Pour activer ou désactiver cet affichage, cliquez sur :

Vue > Balisage de flux

Modification de la taille des icônes d'un flux

Vous pouvez changer la taille des icônes de flux par l'une des méthodes suivantes.

- A l'aide d'un paramètre de propriété de flux
- A l'aide d'un menu contextuel dans le flux
- A l'aide du clavier

Vous pouvez redimensionner la vue entière du flux à une taille comprise entre 8 % et 200 % de la taille d'icône standard.

Pour redimensionner le flux entier (méthode des propriétés du flux)

1. Dans le menu principal, sélectionnez :

Outils > Propriétés du flux > Options > Présentation.

2. Sélectionnez la taille souhaitée dans le menu Taille d'icône.

- 3. Cliquez sur Appliquer pour afficher les résultats.
- 4. Cliquez sur **OK** pour enregistrer les modifications.

Pour redimensionner le flux entier (méthode du menu)

- 1. Cliquez avec le bouton droit de la souris sur l'arrière-plan du flux dans l'espace de travail.
- 2. Sélectionnez l'option Taille d'icône puis la taille souhaitée.

Pour redimensionner le flux entier (méthode du clavier)

- 1. Appuyez sur Ctrl + [-] sur le clavier pour effectuer un zoom arrière et réduire la vue d'une taille.
- 2. Appuyez sur Ctrl + Maj + [+] sur le clavier pour effectuer un zoom avant et agrandir la vue d'une taille.

Cette fonction est particulièrement utile pour obtenir une vue globale d'un flux complexe. Vous pouvez aussi l'utiliser pour réduire le nombre de pages nécessaires à l'impression d'un flux.

Utilisation de la souris dans IBM SPSS Modeler

Dans IBM SPSS Modeler, les utilisations les plus courantes de la souris sont les suivantes :

- Clic simple. Utilisez le bouton droit ou le bouton gauche de la souris pour sélectionner des options dans les menus, ouvrir des menus contextuels, ou accéder à diverses autres commandes et options standard. Cliquez avec la souris et maintenez le bouton de la souris enfoncé pour faire glisser des noeuds.
- **Double-clic.** Cliquez deux fois avec le bouton gauche de la souris pour placer des noeuds dans le canevas de flux et éditer des noeuds existants.
- Clic avec le bouton central. Cliquez avec le bouton central de la souris et faites glisser le curseur pour connecter des noeuds dans le canevas de flux. Double-cliquez avec le bouton central de la souris pour déconnecter un noeud. Si vous ne possédez pas de souris à trois boutons, vous pouvez simuler cette fonction en appuyant sur la touche Alt tout en cliquant avec la souris et en la faisant glisser.

Utilisation des touches de raccourci

Dans IBM SPSS Modeler, de nombreuses opérations de programmation visuelle sont associées à des touches de raccourci. Par exemple, vous pouvez supprimer un noeud en cliquant dessus et en appuyant sur la touche Suppr de votre clavier. De la même façon, vous pouvez enregistrer rapidement un flux en appuyant sur la touche S tout en maintenant la touche Ctrl enfoncée. Les commandes de ce type sont indiquées par Ctrl et une autre touche (par exemple, Ctrl+S).

De nombreuses touches de raccourci sont utilisées dans les opérations Windows standard, telles que Ctrl+X pour couper un élément. Ces raccourcis sont pris en charge dans IBM SPSS Modeler, parallèlement à ceux présentés ci-après, propres à l'application.

Remarque : Dans certains cas, les anciennes touches de raccourci utilisées dans IBM SPSS Modeler sont en conflit avec les touches de raccourci Windows standard. Pour que ces anciennes touches de raccourci fonctionnent, il faut utiliser la touche Alt. Par exemple, Ctrl+Alt+C peut activer ou désactiver la mise en cache.

Touche de raccourci	Fonction
Ctrl+A	Tout sélectionner
Ctrl+X	Couper
Ctrl+N	Permet de créer un flux
Ctrl+O	Permet d'ouvrir un flux
Ctrl+P	Imprimer

Tableau 1. Touches de raccourci prises en charge

Tableau 1.	Touches	de	raccourci	prises	en	charge	(suite)
------------	---------	----	-----------	--------	----	--------	---------

Touche de raccourci	Fonction
Ctrl+C	Copier
Ctrl+V	Coller
Ctrl+Z	Annuler
Ctrl+Q	Permet de sélectionner tous les noeuds situés en aval du noeud sélectionné.
Ctrl+W	Permet de désélectionner tous les noeuds en aval (bascule du raccourci Ctrl+Q)
Ctrl+E	Exécuter à partir d'un noeud sélectionné
Ctrl+S	Permet d'enregistrer le flux en cours
Alt+flèches	Permettent de déplacer les noeuds sélectionnés dans le canevas de flux dans le sens indiqué par la flèche utilisée
Maj+F10	Permet d'ouvrir le menu contextuel du noeud sélectionné

Tableau 2. Touches de raccourci prises en charge pour les anciennes touches d'accès rapide

Touche de raccourci	Fonction				
Ctrl+Alt+D	Permet de dupliquer un noeud				
Ctrl+Alt+L	Permet de charger un noeud				
Ctrl+Alt+R	Permet de renommer un noeud				
Ctrl+Alt+U	Permet de créer un noeud Utilisateur				
Ctrl+Alt+C	Permet d'activer et de désactiver le cache				
Ctrl+Alt+F	Vider le cache				
Ctrl+Alt+X	Développer le super noeud				
Ctrl+Alt+Z	Permet d'effectuer un zoom avant/arrière				
Suppr	Permet de supprimer un noeud ou une connexion				

Impression

Les objets suivants peuvent être imprimés dans IBM SPSS Modeler :

- Diagrammes de flux
- Graphiques
- Tables
- Rapports (à partir du noeud Rapport et des rapports de projet)
- Scripts (à partir des boîtes de dialogue Propriétés du flux, Script autonome ou Script Super noeud)
- Modèles (navigateurs de modèle, onglets de boîte de dialogue avec élément en cours, visualiseurs d'arbres)
- Annotations (à partir de l'onglet Annotations de la sortie)

Pour imprimer un objet :

- Pour imprimer sans afficher d'aperçu, cliquez sur le bouton Imprimer de la barre d'outils.
- Pour définir la mise en page avant d'imprimer, sélectionnez Mise en page dans le menu Fichier.
- Pour afficher un aperçu avant d'imprimer, sélectionnez Aperçu avant impression dans le menu Fichier.
- Pour afficher la boîte de dialogue d'impression standard vous permettant de sélectionner les imprimantes et de définir des options d'aspect, sélectionnez **Imprimer** dans le menu Fichier.

Automatisation d'IBM SPSS Modeler

Etant donné que l'exploration de données avancé peut être complexe et parfois long, IBM SPSS Modeler comprend plusieurs types d'assistance au codage et à l'automatisation.

- **Control Language for Expression Manipulation** (CLEM) est un langage permettant d'analyser et de manipuler les données circulant au sein des flux IBM SPSS Modeler. Les data miners utilisent beaucoup le langage CLEM dans les opérations de flux pour exécuter des tâches aussi simples que le calcul du profit à partir des données de coûts et de revenus, ou aussi complexes que la transformation de données du log Web en un ensemble de champs et d'enregistrements contenant des informations utilisables.
- La génération de scripts est un outil performant pour automatiser les processus dans l'interface utilisateur. Les scripts effectuent des opérations semblables à celles qui peuvent être exécutées à la souris ou au clavier. Vous pouvez également définir une sortie et manipuler des modèles générés.

Chapitre 3. Introduction à la modélisation

Un modèle est un ensemble de règles, de formules, ou d'équations pouvant être utilisées pour prédire un résultat en fonction d'un ensemble de champs ou de variables d'entrée. Par exemple, une institution financière peut utiliser un modèle pour prédire si les emprunteurs représentent un risque important ou peu de risque, en fonction des informations déjà connues sur le passé de ces emprunteurs.

La capacité à prédire un résultat est l'objectif central de l'analyse prédictive, et la compréhension du processus de modélisation est essentielle pour l'utilisation d'IBM SPSS Modeler.



Figure 11. Modèle d'arbre de décision simple

Cet exemple utilise un modèle d'**arbre décision** qui classifie les enregistrements (et prédit une réponse) à l'aide d'une série de règles de décisions, par exemple :

Si revenu = Moyen Et cartes <5 Alors -> 'Bon'

Bien que cet exemple utilise un modèle CHAID (Chi-Squared Automatic Interaction Detection), il est destiné à fournir une introduction générale, et la plupart des concepts s'appliquent globalement aux autres types de modélisation dans IBM SPSS Modeler.

Pour comprendre tous les modèles, vous devez d'abord comprendre les données qu'ils contiennent. Les données de cet exemple contiennent des informations sur les clients d'une banque. Les champs suivants sont utilisés :

Nom du champ	Description
Conditions_crédit	Conditions de crédit : 0=Mauvaises, 1=Bonnes, 9=valeurs manquantes
Age	Age en années
Revenu	Niveau de revenu : 1=Bas, 2=Moyen, 3=Elevé
Cartes_crédit	Nombre de cartes de crédit possédées : 1=Moins de cinq, 2=Cinq ou plus
Education	Niveau d'éducation : 1=Lycée, 2=Université
Prêts_voiture	Nombre de prêts voiture en cours : 1=Aucun ou un, 2=Plus de deux

La banque gère une base de données contenant des informations sur les clients qui ont contracté un prêt, notamment sur le respect de leur engagement de remboursement (conditions de crédit = bonnes) ou le non-respect de leur engagement (conditions de crédit = mauvaises). A l'aide de ces données, la banque peut créer un modèle qui lui permettra de prédire les probabilités de remboursement des futurs emprunteurs.

A partir d'un modèle d'arbre de décision, vous pouvez analyser les caractéristiques de deux groupes de clients et prédire les risques de non-remboursement.

Cet exemple utilise le flux nommé *modelingintro.str*, disponible dans le dossier *Demos* du sous-dossier des *flux*. Le fichier de données est *tree_credit.sav*. Pour plus d'informations, voir la rubrique «Dossier Demos», à la page 4.

Regardons le flux de plus près.

1. Dans le menu principal, sélectionnez les options suivantes :

Fichier > Ouvrir un flux

- Cliquez sur l'icône de la pépite d'or dans la barre d'outils de la boîte de dialogue Ouvrir et choisissez le dossier Demos.
- 3. Double-cliquez sur le dossier des *flux*.
- 4. Double-cliquez sur le fichier modelingintro.str.

Création du flux





Pour construire un flux qui va créer un modèle, vous avez besoin d'au moins trois éléments :

- Un noeud source qui lit les données issues d'une source externe, dans ce cas un fichier de données IBM SPSS Statistics.
- Un noeud source ou type qui spécifie les propriétés des champs, telles que le niveau de mesure (le type de données contenues dans le champ) et le rôle de chaque champ en tant que cible ou entrée dans la modélisation.
- Un noeud de modélisation qui génère un nugget de modèle lors de l'exécution du flux.

Dans cet exemple, nous utilisons un noeud de modélisation CHAID. CHAID, ou Chi-Squared Automatic Interaction Detection, est une méthode de classification qui crée des arbres de décision à l'aide d'un type de statistiques spécifique connu sous le nom de statistiques du khi-deux et qui permet de définir les meilleurs endroit auxquels opérer le découpage dans l'arbre de décision. Si les niveaux de mesure sont spécifiés dans le noeud source, le noeud type distinct peut être éliminé. D'un point de vue fonctionnel, le résultat est le même.

Ce flux comporte également des noeuds Table et Analyse qui seront utilisés pour afficher les résultats de scoring après la création du nugget de modèle et son ajout au flux.

Le noeud Statistics lit les données au format IBM SPSS Statistics à partir du fichier de données *tree_credit.sav*, qui est installé dans le dossier *Demos*. (Une variable spéciale nommée *\$CLEO_DEMOS* est utilisée pour faire référence à ce dossier sous l'installation IBM SPSS Modeler actuelle. Ainsi, le chemin sera toujours valide, quelque soit le dossier d'installation actuel ou la version.)

😵 tree_credit.sav	X
Preview 2 Refresh	0
\$CLEO_DEMOS'tree_credit.sav	
Data Filter Types Annotations	
Import file: \$CLEO_DEMOS'tree_credit.sav	
Variable names: O Read names and labels O Read labels as names	
Values: 🔘 Read data and labels 💿 Read labels as data	
Use field format information to determine storage	
OK	Apply Reset

Figure 13. Lecture des données avec un noeud source Statistics

Le noeud type définit le **niveau de mesure** pour chaque champ. Le niveau de mesure est une catégorie qui indique le type de données du champ. Notre fichier de données source utilise trois niveaux de mesure différents.

Un champ **Continu** (comme le champ *Age*) contient des valeurs numériques continues, alors qu'un champ **Nominal** (comme le champ *Conditions de crédit*) contient deux valeurs distinctes minimum, par exemple *Mauvaises*, *Bonnes*, ou *Pas d'antécédents de crédit*. Un champ **Ordinal** (comme le champ *Niveau de revenu*) décrit les données avec différentes valeurs distinctes ayant un ordre inhérent, dans ce cas *Faible*, *Moyenne* et *Elevée*.

		Y	Ŷ		
	Read Val	lues Clear \	Values	Clear All Valu	les
Field 🗂	Measurement	Values	Missing	Check	Role
Credit rating 🧯	📩 Nominal	Bad,Good	*	None	O Target
👌 Age 💊	🔗 Continuous	[20.00269		None	🔪 Input
Income level 🍙	🚺 Ordinal	High,Low,		None	🔪 Input
Number of 🧯	📩 Nominal	"Less tha		None	🔪 Input
Education	📩 Nominal	"High sch		None	🔪 Input
Carloans	📩 Nominal	"None or		None	hugh in the second

Figure 14. Définition des champs cibles et des champs d'entrées avec le noeud type

Pour chaque champ, le noeud type spécifie également un **rôle**, qui indique le rôle que joue chaque champ dans la modélisation. Le rôle est défini sur *Cible* pour le champ *Conditions de crédit*, qui indique si un client donné a remboursé ou non son prêt. Il s'agit de la **cible**, soit le champ dont vous souhaitez prédire la valeur.

Le rôle est défini sur *Entrée* pour les autres champs. Les champs d'entrée sont quelquefois désignés sous le nom de **prédicteurs**, ou champs dont les valeurs sont utilisées par l'algorithme de modélisation afin de prévoir la valeur du champ cible.

Le noeud de modélisation CHAID génère le modèle.

Sur l'onglet Champs du noeud de modélisation, l'option **Utiliser les rôles prédéfinis** est sélectionnée, ce qui signifie que la cible et les entrées indiquées dans le noeud type seront utilisées. Vous pouvez modifier les rôles de champ à ce stade, mais pour cet exemple, nous les utiliserons tels quels.

1. Cliquez sur l'onglet Options de création.

Cree	ditrating					
CHAID						0
	Objective: Stand	dard model				
Fielde	Duild Outlines	Madel Outlines	0			
TICICIS	Build Options	woder Options	Annotations	-		
O						
O Use	pre <u>d</u> efined roles custom field assi	ianments				
Fields						
Sort:	Jone.				Targets*:	
			1.00			
					Predictors (Inputs)*:	
					Age	
					A Number of credit cards	
			2			
			6	•	💑 Car Ioans	
					Analysis Weight:	
All		1				× & 18 1 4
					26	
OK	📄 🕨 Run	Cancel				Apply Reset

Figure 15. Noeud de modélisation CHAID - Onglet Champs

Plusieurs options sont disponibles ici dans lesquelles vous pouvez spécifier le type de modèle que vous voulez créer.

Nous voulons un tout nouveau modèle, donc nous utiliserons l'option par défaut **Créer un nouveau modèle**.

Nous voulons également un seul modèle d'arbre décision standard sans aucune amélioration, donc nous conserverons l'option d'objectif par défaut **Créer un seul arbre**.

Bien que vous puissiez lancer une session de modélisation interactive qui vous permet d'ajuster le modèle, cet exemple génère simplement un modèle à l'aide du paramètre de mode par défaut **Générer le modèle**.



Figure 16. Noeud de modélisation CHAID - Onglet Options de création

Pour cet exemple, nous avons voulu présenter un arbre relativement simple ; nous limiterons donc la croissance de l'arbre en augmentant le nombre minimum d'observations pour les noeuds parent et enfant.

- 2. Dans l'onglet Options de création, sélectionnez **Règles d'arrêt** dans le panneau de gauche du navigateur.
- 3. Sélectionnez l'option Utiliser la valeur absolue.
- 4. Définissez Enregistrements minimum dans la branche parent sur 400.
- 5. Définissez Enregistrements minimum dans la branche enfant sur 200.


Figure 17. Définition des critères d'arrêt pour la création d'un arbre de décisions

Nous pouvons utiliser toutes les autres options par défaut pour cet exemple, par conséquent, cliquez sur **Exécuter** pour créer le modèle. (Vous pouvez également cliquer avec le bouton droit de la souris sur le noeud et choisir **Exécuter** dans le menu contextuel ou sélectionner le noeud et choisir **Exécuter** dans le menu Outils.)

Navigation dans le modèle

Lorsque l'exécution se termine, le nugget de modèle est ajouté à la palette Modèles dans le coin supérieur droit de la fenêtre de l'application, et est aussi placé dans l'espace de travail du flux avec un lien vers le noeud de modélisation à partir duquel il a été créé. Pour consulter les détails du modèle, cliquez avec le bouton droit de la souris sur le nugget de modèle **Parcourir** (dans la palette des modèles) ou **Modifier** (dans l'espace de travail).

· s	Streams Outputs Models
Add <u>T</u> o Stream	-
Browse	t rating
Rename and Annotate	
🏷 Generate Modeling Node	
Save Model	
Save Model As	
😻 Store Model	
Export PMML	
Add to Project	
× Delete Delete	

Figure 18. Palette Modèles

Dans le cas du nugget CHAID, l'onglet Modèle affiche les détails sous la forme d'un ensemble de règles. Il s'agit essentiellement d'une série de règles pouvant être utilisées pour affecter des enregistrements individuels à des noeuds enfant, en fonction des valeurs des différents champs d'entrée.

A File K Generate of View Preview A	
CIATO	
Model Viewer Summary Settings Annotations	
1 2 All 🖓 📧 🛈	
Income level in ["High"] [Mode: Good]	
Number of credit cards in ["Less than 5"] [Mode: Good] ⇒ Good	
Number of credit cards in ["5 or more"] [Mode: Good] => Good	
Income level in ["Low"] [Mode: Bad] 🖙 Bad	
🖻 Income level in ["Medium"] [Mode: Good]	
— Number of credit cards in ["Less than 5"] [Mode: Good] I Good	
Number of credit cards in ["5 or more"] [Mode: Bad] 🖙 Bad	

Figure 19. Nugget de modèle CHAID, ensemble de règles

Pour chaque noeud de terminal d'arbre de décision, c'est-à-dire ces noeuds Arbre qui ne sont pas plus divisés, une prévision de *Bon* ou *Mauvais*est renvoyée. Dans chaque cas la prédiction est déterminée par le **noeud**, ou par la réponse la plus courante pour les enregistrements qui sont compris dans ce noeud.

A droite de l'ensemble de règles, l'onglet Modèle affiche le graphique d'importance des prédicteurs qui montre l'importance relative de chaque prédicteur dans l'estimation du modèle. Nous pouvons observer que le *niveau de revenu* est le critère plus important dans ce cas et que le seul autre facteur intéressant est le *Nombre de cartes de crédit*.



Figure 20. Graphique de l'importance des prédicteurs

L'onglet Visualiseur dans le nugget de modèle affiche le même modèle sous la forme d'un arbre, avec un noeud à chaque point de décision. Utilisez les commandes du Zoom sur la barre d'outils pour effectuer un zoom avant sur un noeud spécifique ou un zoom arrière pour afficher une plus grande partie de l'arbre.



Figure 21. Onglet Visualiseur dans le nugget de modèle, avec zoom arrière sélectionné

Si l'on regarde la partie supérieure de l'arbre, le premier noeud (Noeud 0) propose un récapitulatif de tous les enregistrements dans le jeu de données. Un peu plus de 40 % des observations de ce jeu de données sont classées comme risquées. Il s'agit d'une proportion élevée. Voyons si l'arbre peut nous donner des informations sur les facteurs responsables.

Nous pouvons observer que la première division se situe au niveau du *Niveau de revenu*. Les enregistrements dans lesquels le niveau de revenu se trouve dans la catégorie *Low* (Faible) sont affectés au Noeud 2 et il n'est pas surprenant de voir que cette catégorie contient le plus fort pourcentage de non-reboursement de prêts. Il est évident qu'accorder un prêt aux clients de cette catégorie présente un risque élevé.

Cependant, 16% des clients de cette catégorie ont, en réalité, *remboursé* leur prêt. Par conséquent, cette prévision n'est pas toujours exacte. Aucun modèle ne peut réellement prédire toutes les réponses, mais un bon modèle doit vous permettre de prédire la réponse *la plus problable* pour chaque enregistrement, sur la base des données disponibles.

De la même façon, si l'on observe les clients avec un revenu élevé (Noeud 1), on s'aperçoit que la grande majorité (89 %) présente un risque peu élevé. Mais plus de 1 clients sur 10 n'a pas remboursé son prêt. Est-il possible d'affiner nos critères de prêt pour diminuer le risque ?

Veuillez noter que le modèle a divisé ces clients en deux sous-catégories (noeuds 4 et 5), en fonction du nombre de cartes de crédit possédées. Pour les clients à revenu élevé, si nous prêtons uniquement à ceux possédant moins de 5 cartes de crédit, nous pouvons faire passer notre taux de succès de 89% à 97%, soit un résultat encore plus satisfaisant.



Figure 22. Vue sous forme d'arbre des clients à revenu élevé

Mais qu'en est-il des clients appartenant à la catégorie Revenu moyen (Noeud 3) ? Ils sont encore plus fortement divisés entre les évaluations Good (Bon) et Bad (Mauvais).

De nouveau, les sous-catégories (Noeuds 6 et 7 dans ce cas) peuvent nous aider. Cette fois, prêter uniquement aux clients avec des revenus moyens et possédant moins de 5 cartes de crédit fait passer le pourcentage de conditions Bonnes de 58% à 85%, soit une augmentation importante.



Figure 23. Vue sous forme d'arbre des clients à revenu moyen

Nous avons appris que chaque enregistrement contenu dans ce modèle sera attribué à un noeud spécifique et recevra une prévision *Bonne* ou *Mauvaise* en fonction des réponses les plus courantes de ce noeud.

Ce processus consistant à affecter des prédictions à des enregistrements individuels s'appelle le **scoring** (ou évaluation). En effectuant le scoring des mêmes enregistrements utilisés pour estimer le modèle, il est possible d'évaluer son exactitude sur les données d'apprentissage, données dont nous connaissons le résultat. Examinons comment effectuer cette opération.

Evaluation du modèle

Nous avons parcouru le modèle pour comprendre le fonctionnement du scoring. Mais pour évaluer son *exactitude*, nous devons déterminer le score de certains enregistrements et comparer les réponses prédites par le modèle aux résultats réels. Nous allons déterminer le score des mêmes enregistrements qui ont été utilisés pour estimer le modèle, ce qui nous permet de comparer les réponses observées et les réponses prédites.



Figure 24. Lier le nugget de modèle au noeuds de sortie pour l'évaluation du modèle

1. Pour voir les scores ou les prédictions, attachez le noeud Table au nugget de modèle, double-cliquez sur le noeud Table et cliquez sur **Exécutez**.

La table affiche les scores prédits dans un champ nommé *R-Credit rating*, qui a été créé par le modèle. Nous pouvons comparer ces valeurs au champ *Conditions de crédit* d'origine qui contient les réponses réelles.

Par convention, les noms des champs générés au cours du scoring sont basés sur le champ cible, mais avec un préfixe standard. Les préfixes G et GE sont générés par le Modèle linéaire généralisé, R est le préfixe utilisé pour la prévision générée par le modèle CHAID dans ce cas, RC est utilisé pour les valeurs de confiance, X est généralement utilisé en utilisant un ensemble, et XR, XS et XF sont utilisés respectivement comme préfixes lorsque le champ cible est un champ de type Continu, Catégoriel, Ensemble ou Indicateur. Les types de modèles différents utilisent des ensembles de préfixes distincts. Une **valeur de confiance** est la propre estimation du modèle, sur une échelle de 0,0 à 1,0, de l'exactitude de chaque valeur prédite.

				0
			12	
Number of credit cards	Education	Car loans	\$R-Credit rating	\$RC-Credit rating
5 or more	College	More than 2	Bad	0.560
5 or more	College	More than 2	Bad	0.560
5 or more	High school	More than 2	Bad	0.832
5 or more	College	None or 1	Bad	0.832
5 or more	College	More than 2	Bad	0.560
5 or more	College	More than 2	Bad	0.560
5 or more	College	More than 2	Bad	0.560
5 or more	High school	More than 2	Bad	0.832
5 or more	High school	More than 2	Bad	0.832
5 or more	College	More than 2	Bad	0.560
5 or more	College	More than 2	Bad	0.832
5 or more	High school	More than 2	Bad	0.832
5 or more	High school	More than 2	Bad	0.560
5 or more	College	None or 1	Bad	0.832
5 or more	High school	More than 2	Bad	0.832
5 or more	College	More than 2	Bad	0.832
5 or more	College	More than 2	Bad	0.832
5 or more	College	More than 2	Bad	0.560
5 or more	College	More than 2	Bad	0.560
5 or more	College	More than 2	Good	0.827

Figure 25. Table affichant les scores générés et les valeurs de confiance

Comme prévu, la valeur prédite correspond aux réponses réelles pour de nombreux enregistrements mais pas pour tous. La raison à cela est que chaque noeud terminal CHAID comporte un ensemble de réponses. La prédiction correspond à la réponse la *plus courante*, mais elle sera fausse pour toutes les autres réponses de ce noeud. (Pensez à la minorité de 16% de clients à faible revenu qui ont remboursé leur prêt).

Pour éviter ceci, nous pouvons continuer à diviser l'arbre en branches de plus en plus petites, jusqu'à ce que chaque noeud soit pur à 100%, autrement dit qu'il ne comporte que des *Good* (Bonne) ou *Bad* (Mauvaise) sans réponses mixtes. Mais un tel modèle serait extrêmement compliqué et serait probablement difficile à étendre à d'autres jeux de données.

Pour connaître précisément le nombre de prévisions correctes, nous pouvons lire la table et compter le nombre d'enregistrements où la valeur du champ prédit *R-Credit rating* correspond à la valeur des *Conditions de crédit*. Heureusement, il y a beaucoup plus simple : nous pouvons utiliser le noeud Analyse, qui effectue automatiquement cette opération.

- 2. Connectez le nugget de modèle au noeud Analyse.
- 3. Double-cliquez sur le noeud Analyse, puis cliquez sur Exécuter.



Figure 26. Ajout d'un noeud Analyse

L'analyse montre que pour 1899 enregistrements sur 2464 - un peu plus de 77% - la valeur prédite par le modèle correspondait à la réponse réelle.

🔍 Analysi	s of [Credit	rating]		
📦 <u>F</u> ile	è Edit 🛛 🚺		}	0 ×
Analysis	Annotations			
Collaps	se All 🗣 E	Expand All		
-Results	for output field	Credit rating	3	
<u>⊟</u> -Con	nparing \$R-Cre	dit rating wit	h Credit rating	
	Correct	1,899	77.07%	
	Wrong	565	22.93%	
	Total	2,464		
				ОК

Figure 27. Résultats d'analyse comparant les réponses observées et les réponses prédites

Ce résultat est limité parce que les enregistrements auxquels un score est donné sont les mêmes que ceux utilisés pour évaluer le modèle. Dans la réalité, vous pourriez utiliser un noeud Partitionner pour diviser les données en échantillons distincts pour l'apprentissage et l'évaluation.

L'utilisation d'un échantillon de partition pour la génération du modèle et d'un autre échantillon pour le tester vous permet d'avoir une bien meilleure indication de la manière dont il peut s'étendre à d'autres jeux de données.

Le noeud Analyse nous permet de tester le modèle sur les enregistrements pour lesquels nous connaissons déjà le résultat réel. L'étape suivante illustre la façon dont nous pouvons utiliser le modèle pour évaluer les enregistrements dont nous ne connaissons pas le résultat. Par exemple, cela peut comprendre les gens qui ne sont pas des clients de la banque, mais qui sont des cibles potentielles pour un publipostage promotionnel.

Scoring des enregistrements

Auparavant, nous avons évalué les mêmes enregistrements utilisés pour estimer le modèle afin de connaître l'exactitude du modèle. A présent, nous allons voir comme évaluer un ensemble d'enregistrements différent de ceux utilisés pour créer le modèle. Il s'agit de l'objectif de la modélisation avec un champ cible : étudier les enregistrements pour lesquels vous connaissez le résultat pour identifier des schémas qui vous permettront de prédire les résultats que vous ne connaissez pas encore.



Figure 28. Association de nouvelles données pour le scoring

Vous pouvez mettre à jour le noeud source Statistics pour qu'il pointe vers un fichier de données différent ou vous pouvez ajouter un nouveau noeud source qui lit dans les données que vous voulez évaluer. Dans les deux méthodes, le nouveau jeu de données doit contenir les mêmes champs d'entrée utilisés par le modèle (*Age, Niveau de revenu, Education,* etc.) mais pas le champ cible *Conditions de crédit*.

Vous pouvez également ajouter le nugget de modèle à tout flux contenant les champs d'entrée attendus. Qu'il soit lu à partir d'un fichier ou d'une base de données, le type de source n'importe pas du moment que les noms et les types des champs correspondent à ceux utilisés par le modèle.

Vous pouvez également enregistrer le nugget de modèle en tant que fichier distinct, exporter le modèle au format PMML pour une utilisation avec d'autres applications qui prennent en charge ce format ou stocker le modèle dans un répertoire IBM SPSS Collaboration and Deployment Services, ce qui permet le déploiement, le scoring et la gestion des modèles à l'échelle de l'entreprise.

Quelque soit l'infrastructure utilisée, le modèle proprement dit fonctionne de la même manière.

Récapitulatif

Cet exemple décrit la procédure standard de création, d'évaluation et de scoring d'un modèle.

- Le noeud de modélisation estime le modèle en étudiant les enregistrements pour lesquels le résultat est connu et crée un nugget de modèle. On parle parfois d'apprentissage du modèle.
- Le nugget de modèle peut être ajouté à n'importe quel flux contenant les champs attendus pour évaluer les enregistrements. En effectuant le scoring des enregistrements pour lesquels vous connaissez déjà le résultat (les clients existants par exemple), vous pouvez évaluer la performance du modèle.

- Une fois que vous êtes satisfait de la performance du modèle, vous pouvez effectuer un scoring de nouvelles données (des clients potentiels par exemple) pour prédire leur réponse.
- Les données utilisées pour l'apprentissage ou l'estimation du modèle peuvent être appelées données analytiques ou historiques; les données de scoring peuvent également être appelées données opérationnelles.

Chapitre 4. Modélisation automatisée d'une cible indicateur

Modélisation de la réponse client (Discriminant automatique)

Le noeud Discriminant automatique vous permet de créer et de comparer automatiquement différents modèles pour les cibles indicateur (comme la probabilité selon laquelle un client donné est susceptible ou non de rembourser une échéance de prêt ou de répondre à une offre spécifique) ou les cibles (d'ensemble) nominales . Dans cet exemple, nous allons rechercher un résultat indicateur (oui ou non). Dans un flux relativement simple, le noeud génère et classe un ensemble de modèles candidats, choisit les meilleurs et les combine en un modèle (combiné) agrégé unique. Cette approche conjugue la facilité de l'automatisation aux avantages de combiner plusieurs modèles ce qui permet généralement des prédictions plus précises que celles de tout autre modèle.

Cet exemple repose sur une société fictive qui souhaite obtenir des résultats plus rentables en présentant à chaque client une offre adaptée.

Cette approche souligne les avantages de l'automatisation. Pour un exemple similaire qui utilise une cible continue (plage numérique), voir Valeurs des propriétés (Numérisation automatique).



Figure 29. Flux d'échantillons Discriminant automatique

Cet exemple utilise le flux *pm_binaryclassifier.str*, installé dans le dossier Démo dans le répertoire des *flux*. Le fichier de données est *pm_customer_train1.sav*. Pour plus d'informations, voir la rubrique «Données d'historique».

Données d'historique

Le fichier *pm_customer_train1.sav* comporte des données d'historique suivant les offres faites à des clients spécifiques au cours de campagnes passées, comme l'indique la valeur du champ *campaign*. Le plus grand nombre d'enregistrements se trouve dans la campagne *Premium account*.

Les valeurs du champ *campaign* sont en fait codées comme des entiers dans les données (par exemple 2 = *Premium account*). Plus tard, vous définirez les libellés de ces valeurs qui vous permettront d'obtenir des résultats plus probants.

違 File	📄 <u>E</u> dit (<u>G</u> enerate	. 🔒					0
Table ,	Annotations							
	customer_id	campaign	response	response_date	purchase	purchase_date	product_id	Rowid
1	7	2	0	\$null\$	0	\$null\$	\$null\$	1
2	13	2	0	\$null\$	0	\$null\$	\$null\$	2
3	15	2	0	\$null\$	0	\$null\$	\$null\$	3
4	16	2	1	2006-07-05 00:00:00	0	\$null\$	183	761
5	23	2	0	\$null\$	0	\$null\$	\$null\$	4
6	24	2	0	\$null\$	0	\$null\$	\$null\$	5
7	30	2	0	\$null\$	0	\$null\$	\$null\$	6
8	30	3	0	\$null\$	0	\$null\$	\$null\$	7
9	33	2	0	\$null\$	0	\$null\$	\$null\$	8
10	42	3	0	\$null\$	0	\$null\$	\$null\$	9
11	42	2	0	\$null\$	0	\$null\$	\$null\$	10
12	52	2	0	\$null\$	0	\$null\$	\$null\$	11
13	57	2	0	\$null\$	0	\$null\$	\$null\$	12
14	63	2	1	2006-07-14 00:00:00	0	\$null\$	183	1501
15	74	2	0	\$null\$	0	\$null\$	\$null\$	13
16	74	3	0	\$null\$	0	\$null\$	\$null\$	14
17	75	2	0	\$null\$	0	\$null\$	\$null\$	15
18	82	2	0	\$null\$	0	\$null\$	\$null\$	16
19	89	3	0	\$null\$	0	\$null\$	\$null\$	17
20	89	2	0	\$null\$	0	\$null\$	\$null\$	18
				And and a first		Accession Con-	de la constanción de	1

Figure 30. Données sur les anciennes promotions

Ce fichier comprend également un champ *réponse* qui indique si l'offre a été acceptée (0 = *non*, et 1 = *oui*). Il s'agit du **champ cible**, ou la valeur, que vous souhaitez prédire. Plusieurs champs contenant des informations démographiques et financières sur chaque client ont également été ajoutés. Ces champs peuvent permettre de créer ou de "former" un modèle qui prédit les taux de réponse des individus ou des groupes en fonction de caractéristiques telles que le revenu, l'âge ou le nombre de transactions mensuelles.

Création du flux

Ajoutez un noeud source Statistiques qui pointe sur *pm_customer_train1.sav*, dans le dossier *Demos* du répertoire d'installation d'IBM SPSS Modeler. (Vous pouvez saisir \$CLE0_DEMOS/ dans le chemin d'accès comme raccourci permettant de référencer ce dossier. Veuillez noter qu'une barre oblique (/) plutôt qu'une barre oblique inverse (\) doit être utilisée dans le chemin d'accès, comme indiqué.)



Figure 31. Lecture de données

2. Ajoutez un noeud type, puis sélectionnez *Réponse* en tant que champ cible (Rôle = **Cible**). Paramétrez l'option Mesure de ce champ sur **Indicateur**.

Type	Annotations				× • •
~	🕨 🛛 🌔 Read Va	lues Clear	Values	Clear All Va	alues
Field -	Measurement	Values	Missing	Check	Role
父 customer_id 🤞	🔗 Continuous	[7,116993]	-	None	🛇 None 🔺
🚫 campaign 🧯	Nominal	1,2,3,4		None	🔪 Input
🚫 response 🛛 🖁	🖌 Flag	1/0		None	🔘 Target
💼 response 🖌	🔗 Continuous	[2006-04		None	O None
🚫 purchase 🛛 🖌	🔗 Continuous	[0,1]		None	○ None
🔁 purchase 🖌	🔗 Continuous	[2006-04		None	○ None
🚫 product_id 🛛 🖌	🔗 Continuous	[183,421]		None	○ None
📿 Rowid 🛛 🖌	🔗 Continuous	[1,19599]		None	○ None
🛆 ana 🛆	🖉 Continuous	110 OE1		None	🔪 Innut 🔼
 View current fi OK Cancel 	elds 🔘 View unu	sed field setting	gs		Apply Reset

Figure 32. Configuration du niveau de mesure et du rôle

- **3**. Définissez l'option role sur **None** (aucun) pour les champs suivants : *customer_id, campaign, response_date, purchase, purchase_date, product_id, Rowid* et X_*random*. Ces champs seront ignorés lors de la création du modèle.
- 4. Cliquez sur le bouton Lire les valeurs dans le noeud type pour vérifier que les valeurs sont instanciées.

Comme nous l'avons vu auparavant, nos données source contiennent des informations sur quatre campagnes différentes, chacune visant un type de compte client différent. Ces campagnes sont codées comme entiers dans les données, et pour se rappeler plus facilement quel type de compte

chaque entier représente, définissons les libellés de chacun d'eux.

Type Types Format Annotations				0	
Read Va	alues Clear	Values	Clear All Val	lues	
Field - Measurement	Values	Missing	Check	Role	
🔗 customer id 🔗 Continuous	[7,116993]	h	lone	♦ None	4
🚫 campaign 🛛 🂑 Nominal	<curr td="" 🔽<=""><td>h</td><td>lone</td><td>🔪 Input</td><td></td></curr>	h	lone	🔪 Input	
🚫 response 🖁 Flag	<read></read>	ľ	lone	O Target	
response 🔗 Continuous	<read +=""></read>		lone	O None	
🐼 purchase 🔗 Continuous	«Pass»		lone	O None	
purchase 🔗 Continuous	«Current»	P	lone	O None	
🐼 product_id 🔗 Continuous	Specify .	P	lone	O None	
🐼 Rowid 🧳 Continuous	TT, 19599, 1	P	lone	○ None	
🛆 ana 🛛 🔊 Continuous	190 011	h	lone	N Innut	-
View current fields View unu OK Cancel	ised field setting	ļS	[eset

Figure 33. Choix de spécification des valeurs d'un champ

- 5. Sur la ligne du champ campaign, cliquez sur l'entrée dans la colonne Valeurs.
- 6. Sélectionnez Spécifier dans la liste déroulante.

		A	
Measurement:	ominal Storage:	Model Field	
Values:	Read from data	Pass	
	Specify values and labels		
	- · ·	Labela	
	Values A	Labels	
	2	Premium account	
	3	Gold account	
	4	Platinum account	
		•	
Check values:			
Check values:	None Ks		1
Check values:	None Ks Missing values		
Check values:	None ks Missing values		
Check values:	None Missing values		×
Check values:	ks Missing values Range	to:	×
Check values:	ks Missing values Range Null White space	to:	×
Check values:	ks Missing values Range Null White space	to:	

Figure 34. Définition de libellés pour les valeurs de champ

7. Dans la colonne **Libellés**, saisissez les libellés comme indiqué pour chacune des quatre valeurs du champ **campaign**.

8. Cliquez sur OK.

Vous pouvez maintenant afficher les libellés dans les fenêtres de sortie plutôt que les entiers.

🔰 <u>F</u> ile	📄 <u>E</u> dit 🛛 💐) <u>G</u> enerate 🛛 🚺		a aa			ୄୄୄୄ
Fable	Annotations						
	customer_id	campaign	response	response_date	purchase	purchase_date	product_id
1	7	Premium account	0	\$null\$	0	\$null\$	\$null\$
2	13	Premium account	0	\$null\$	0	\$null\$	\$null\$
3	15	Premium account	0	\$null\$	0	\$null\$	\$null\$
4	16	Premium account	1	2006-07-05 00:00:00	0	\$null\$	183
5	23	Premium account	0	\$null\$	0	\$null\$	\$null\$
6	24	Premium account	0	\$null\$	0	\$null\$	\$null\$
7	30	Premium account	0	\$null\$	0	\$null\$	\$null\$
8	30	Gold account	0	\$null\$	0	\$null\$	\$null\$
9	33	Premium account	0	\$null\$	0	\$null\$	\$null\$
10	42	Gold account	0	\$null\$	0	\$null\$	\$null\$
11	42	Premium account	0	\$null\$	0	\$null\$	\$null\$
12	52	Premium account	0	\$null\$	0	\$null\$	\$null\$
13	57	Premium account	0	\$null\$	0	\$null\$	\$null\$
14	63	Premium account	1	2006-07-14 00:00:00	0	\$null\$	183
15	74	Premium account	0	\$null\$	0	\$null\$	\$null\$
16	74	Gold account	0	\$null\$	0	\$null\$	\$null\$
17	75	Premium account	0	\$null\$	0	\$null\$	\$null\$
18	82	Premium account	0	\$null\$	0	\$null\$	\$null\$
19	89	Gold account	0	\$null\$	0	\$null\$	\$null\$
20	89	Premium account	0	\$null\$	0	\$null\$	\$null\$
	4			derinconta di			•

Figure 35. Affichage de libellés de valeur de champ

- 9. Reliez un noeud Table au noeud type.
- 10. Ouvrez le noeud Table, puis cliquez sur Exécuter.
- 11. Dans la fenêtre de sortie, cliquez sur le bouton de la barre d'outils **Afficher les libellés de champ et de valeur** pour afficher les libellés.
- 12. Cliquez sur OK pour fermer la fenêtre de sortie.

Les données incluent des informations sur quatre campagnes différentes, mais vous vous concentrerez sur l'analyse d'une seule campagne à la fois. Comme le plus grand nombre d'enregistrements se trouve dans la campagne de compte Premium (codée *campaign=2* dans les données), vous pouvez utiliser un noeud Sélectionner pour n'inclure que ces enregistrements dans le flux.

Select		\mathbf{X}
?>	Preview	0
Settings	Annotations	
Mode:	🔘 Include 🔘 Discard	
Condition:	campaign = 2	
OK Car	ncel	Apply

Figure 36. Sélection d'enregistrements pour une seule campagne

Génération et comparaison de modèles

- 1. Liez un noeud Discriminant automatique et sélectionnez **Exactitude globale** comme système métrique utilisé pour classer les modèles.
- 2. Définissez le **Nombre de modèles à utiliser** sur 3. Cela signifie que les trois meilleurs modèles seront créés lorsque vous exécuterez le noeud.

💟 response	×
	0
Estimated number of models to be executed: 9	
Fields Model Expert Discard Settings Annotations	
Model name: O Auto O Custom	
☑ Use partitioned data	
👿 Build model for each split	
Rank models by: Overall accuracy 💎	
Rank models using: O Training partition O Test partition	
Number of models to use:	
Calculate predictor importance	
Profit Criteria (valid only for flag targets)	
Costs:	-
Revenue:	
Weight: Fixed 1.0 Variable	-
Lift Criteria (valid only for flag targets)	
Percentile to use for lift calculation: 30	
OK Frun Cancel	Apply Reset

Figure 37. Noeud Discriminant automatique - Onglet Modèle

Dans l'onglet Expert, vous pouvez choisir jusqu'à 11 algorithmes de modèle différents.

3. Désélectionnez les types de modèle **Discriminant** et **SVM**. (Ces modèles prennent plus longtemps à se former à partir de ces données et les désélectionner accélérera l'exemple. Mais si patienter ne vous dérange pas, n'hésitez pas à les laisser sélectionnés).

Comme vous avez défini le **Nombre de modèles à utiliser** sur 3 dans l'onglet Modèle, le noeud calcule l'exactitude des neuf algorithmes restants et crée un nugget de modèle unique contenant les trois plus précis.

	Estimated	d number	of models to	be executed: 9	
Fields I	Model E:	xpert Di	scard Set	tings Annotations	
lodels us	ed:	Modelts	ma	Model peremeters	No of models
		C Sto	C5	Default	1
	v	K	Logistic r	Default	1
[v	3	Decision	Default	1
	v	*	Bayesian	Default	1
I			Discrimin	Default	1
	~	14	KNN Alg	Default	1
[1100	SVM	Default	1
I	~	^C _{RT}	C&R Tree	Default	1
	-	OWEST	Quest Tr	Default	1
	V	CHAID	CHAID Tree	Default	1
Restric	ct maximu	n time spe	ent building :	a single model to	15 🖨 minutes

Figure 38. Noeud Discriminant automatique - Onglet Expert

4. Dans l'onglet Paramètres, pour la méthode d'ensemble, sélectionnez **Vote pondéré par la confiance**. Cela détermine la façon dont un score agrégé unique est produit pour chaque enregistrement.

Avec le vote simple, si deux modèles sur trois prédisent *oui*, alors*oui* l'emporte par un vote de 2 contre 1. Dans le cas de vote pondéré par la fiabilité, les votes sont pondérés en fonction de la valeur de confiance de chaque prévision. Par conséquent, si un modèle prévoit *non* avec un niveau de confiance plus élevé que les deux prévisions *oui* combinées, alors *non* l'emporte.

	Estin	nated nu	mber of r	nodels to	be executed: §	0 – [
Fields	Model	Expert	Discard	Settings	Annotations	
	emple m oting is ti Randon Raw pr	etnod: ed, select n selectio opensity	t value usir	e-weighted ng: est confide	nce	

Figure 39. Noeud Discriminant automatique - Onglet Paramètres

5. Cliquez sur Exécuter.

Après quelques minutes, le nugget de modèle généré est créé et placé sur l'espace de travail et dans la palette Modèles en haut à droite de la fenêtre. Vous pouvez parcourir le nugget de modèle ou l'enregistrer ou le déployer de plusieurs façons.

Ouvrez le nugget de modèle ; il répertorie les détails concernant chacun des modèles créés au cours de l'exécution. (En situation réelle, lorsque des centaines de modèles peuvent être créés à partir d'un grand nombre de données, cette opération peut prendre plusieurs heures.) Voir figure 29, à la page 39.

Si vous souhaitez analyser plus en détail l'un des modèles individuels, vous pouvez double-cliquer sur un nugget de modèles dans la colonne **Modèle** pour la faire défiler et parcourir les résultats du modèle individuel ; à partir de là, vous pouvez générer des noeuds de modélisation, des nuggets de modèle ou des graphiques d'évaluation. Dans la colonne **Graphique**, vous pouvez double-cliquer sur une miniature pour générer un graphique en grandeur nature.

😡 res	😭 response 🛛 🛛 🔀								
	👔 File 🖏 Generate 🖋 View 🕞 Preview) 🚯 🔹 🕞								
Model	Model Graph Summary Settings Annotations								
Sort by	c Overall acc	curacy 🔻 🛇 Ascendi	ng 🔘 Descendi	ing	X Delet	e Unused Mode	ls View:	Training set *	
Use?	Graph	Model	Build Time (mins)	Max Profit	Max Profit Occurs in (%)	Lift{Top 30%}	Overall Accuracy ∇	No. Fields Used	Area Under Curve
		C5 1	<1	4,906.667	8	2.203	92.861	10	0.777
		C&R Tree 1	3	4,602.692	9	2.778	92.365	8	0.924
	CHAID Tree 1 3 4,145.668 8 2.851 91.706 4 0.927								
ОК	Cancel							Ar	ply <u>R</u> eset

Figure 40. Résultats du Discriminant automatique

Par défaut, les modèles sont classés en fonction de l'exactitude globale, cette mesure ayant été sélectionnée dans l'onglet Modèle du noeud Discriminant automatique. Le modèle C51 se classe en meilleure position selon cette mesure, mais les modèles Arbre C&R et CHAID sont presque aussi précis.

Vous pouvez effectuer le tri sur une autre colonne en cliquant sur l'en-tête de cette colonne ou vous pouvez choisir la mesure désirée dans la liste déroulante **Trier par** de la barre d'outils.

En fonction de ses résultats, vous pouvez décider d'utiliser les trois modèles les plus précis. En combinant les prévisions à partir de plusieurs modèles, il est possible d'éviter les limitations dans les modèles individuels. Ce qui entraîne une plus grande exactitude globale.

Dans la colonne Utiliser ?, sélectionnez les modèles C51, Arbre C&R et CHAID.

Liez un noeud Analyse (palette Sortie) après le nugget de modèle. Cliquez avec le bouton droit de la souris sur le noeud Analyse et sélectionnez **Exécuter** pour exécuter le flux.

Le score agrégé généré par le modèle combiné est affiché dans un champ nommé *\$XF-response*. Lorsque les valeurs prédites sont mesurées en fonction des données d'apprentissage, elles correspondent à la réponse réelle (comme enregistrées dans le champ *réponse* d'origine) avec une exactitude globale de 92,82 %.

Bien que ce modèle ne soit pas aussi précis que le meilleur des trois modèles individuels (92,86 % pour C51), la différence est trop minime pour être significative. Généralement, un modèle combiné sera plus performant lorsqu'il sera appliqué à des jeux de données autres que les données d'apprentissage.

🕙 Analysis of [response] 📃 🗖 🔀						
📦 <u>F</u> ile	🖹 Edit 🛛 🛃		0 ×			
Analysis	Analysis Annotations					
8 Collaps	se All 🤷 E	xpand All				
Results	for output field r	response				
📄 Con	nparing \$XF-res	ponse with r	esponse			
	Correct	12,534	92.82%			
	Wrong	970	7.18%			
	Total	13,504				
				ОК		

Figure 41. Analyse des trois modèles combinés

Récapitulatif

Pour résumer, vous avez utilisé le noeud Discriminant automatique pour comparer plusieurs modèles différents, vous avez utilisé les trois modèles les plus précis et vous les avez ajoutés au flux dans un nugget de modèle Discriminant automatique combiné.

- Concernant l'exactitude globale, les modèles C51, Arbre C&R et CHAID sont plus performants avec les données d'apprentissage.
- Le modèle combiné a presque été aussi performant que le meilleur des modèles individuels et peut être aussi efficace lorsqu'il est appliqué à d'autres jeux de données. Si votre objectif est d'automatiser autant que possible le processus, cette approche vous permet d'obtenir un modèle fiable dans la plupart des circonstances sans avoir à creuser trop dans les spécificités des modèles.

Chapitre 5. Modélisation automatisée d'une cible continue

Valeurs de propriété (Numérisation automatique)

Le noeud Numérisation automatique vous permet de créer et de comparer automatiquement différents modèles pour des résultats continus (intervalle numérique), tels que la prévision de la valeur imposable d'une propriété. Avec un seul noeud, vous pouvez estimer et comparer un ensemble de modèles candidats et générer un sous-ensemble de modèles pour des analyses ultérieures. Ce noeud fonctionne de la même manière que le noeud Discriminant automatique mais pour les cibles continues plutôt que pour les cibles indicateurs ou les cibles nominales.

Le noeud combine le meilleur des modèles candidats dans un nugget de modèle agrégé (d'ensemble) unique. Cette approche conjugue la facilité de l'automatisation aux avantages de combiner plusieurs modèles ce qui permet généralement des prédictions plus précises que celles de tout autre modèle.

Cet exemple se concentre sur un responsable de municipalité fictif qui ajuste et estime les taxes foncières. Pour obtenir une plus grande précision, il va construire un modèle qui prédit les valeurs immobilières en fonction du type de bâtiment, du voisinage, de la taille et d'autres facteurs connus.



Figure 42. Flux d'échantillons Numérisation automatique

Cet exemple utilise le flux *property_values_numericpredictor.str*, installé dans le dossier Démos dans le répertoire des *flux*. Le fichier de données utilisé est *property_values_train.sav*. Pour plus d'informations, voir la rubrique «Dossier Demos», à la page 4.

Données d'apprentissage

Le fichier de données comprend un champ nommé *taxable_value*, qui est le **champ cible**, ou la valeur à prédire. Les autres champs contiennent des informations telles que le voisinage, le type de bâtiment et le volume intérieur et peuvent être utilisés comme prédicteurs.

Nom du champ	Libellé
property_id	ID propriété
neighborhood	Zone à l'intérieur de la ville
building_type	Type de bâtiment
year_built	Année de construction
volume_interior	Volume intérieur
volume_other	Volume du garage et des bâtiments supplémentaires
lot_size	Taille du lot

Nom du champ	Libellé
taxable_value	Valeur imposable

Le dossier Demos contient également un fichier de données de scoring nommé *property_values_score.sav*. Ce fichier contient les mêmes champs mais sans le champ *taxable_value*. Après la formation des modèles à l'aide des jeux de données où la valeur imposable est connue, vous pouvez évaluer des enregistrements où cette valeur ne l'est pas.

Création du flux

1. Ajoutez un noeud source Statistiques qui pointe sur *property_values_train.sav*, dans le dossier *Demos* du répertoire d'installation d'IBM SPSS Modeler. (Vous pouvez saisir \$CLE0_DEMOS/ dans le chemin d'accès comme raccourci permettant de référencer ce dossier. Veuillez noter qu'une barre oblique (/) plutôt qu'une barre oblique inverse (\) doit être utilisée dans le chemin d'accès, comme indiqué.)



Figure 43. Lecture de données

2. Ajoutez un noeud type, puis sélectionnez *taxable_value* en tant que champ cible (Rôle = **Cible**). Le rôle doit être défini sur **Entrée** pour tous les autres champs, indiquant ainsi qu'ils seront utilisés comme prédicteurs.

Types Format	Annotations			Class All Make	
Eield -	Maggurament	Values	Missing	Check	Role
		12 24 44 91	wissing	Mana	
v property_ia		[2,21410]	<u> </u>	None	a input
A neighborhood	nominal	Bioemenp	<u> </u>	None	
A building_type	nominal	"2-onder		None	Input
📿 year_built 🛛 .	🖉 Continuous	[1870,1992]	*	None	🔪 Input
父 volume_inte ,	🖉 Continuous	[138,1901]	*	None	🔪 Input
🚫 volume_other .	🔗 Continuous	[0,496]		None	🔪 Input
🚫 lot_size 🛛 .	🔗 Continuous	[55,1310]	*	None	🔪 Input
🔆 taxable_value .	Continuous	[40000,66	*	None	O Target
View current fi OK Cancel	ields 🔘 View unuse	d field settings		4	Apply <u>R</u> eset

Figure 44. Définition du champ cible

- **3**. Liez un noeud Numérisation automatique et sélectionnez **Corrélation** comme mesure utilisée pour classer les modèles.
- 4. Définissez le **Nombre de modèles à utiliser** sur 3. Cela signifie que les trois meilleurs modèles seront créés lorsque vous exécuterez le noeud.



Figure 45. Noeud Numérisation automatique - Onglet Modèle

5. Dans l'onglet Expert, laissez les paramètres par défaut ; le noeud estime un modèle unique pour chaque algorithme, pour un total de sept modèles. (Vous pouvez également modifier ces paramètres pour comparer plusieurs variantes pour chaque type de modèle.)

Comme vous avez défini le **Nombre de modèles à utiliser** sur 3 dans l'onglet Modèle, le noeud calcule l'exactitude des sept algorithmes restants et crée un nugget de modèle simple contenant les trois plus précis.

Estimated	I number of models to be ex	ecuted: 7			
Models used:					
User	Regression	Default	1		
	Generalized	Default	1		
	KNN Algorithm	Default	1		
	SVM	Default	1		
	KT C&R Tree	Default	1		
	CHAID Tree	Default	1		
	Neural Net	Default	1		

Figure 46. Noeud Numérisation automatique - Onglet Expert

6. Dans l'onglet Paramètres, laissez les paramètres par défaut tels quels. Parce qu'il s'agit d'une cible continue, le score d'ensemble est généré en effectuant la moyenne de ces scores pour les modèles individuels.

😡 taxa	ble_va	lue					×
							0
¥\$	Estim	nated nu	mber of n	nodels to be	executed: 6	3	
Fields	Model	Expert	Settings	Annotations			
Ensem	ble Settir	ngs					
The e	nsemble	scores fo	or a contin	uous target wil	l be generat	ed by averag	ing.
Ca	loulate e	tandard e	rror				
	iouiuto o	tanaara e					
ОК	🕨 Run	Cano	el)				Apply Reset

Figure 47. Noeud Numérisation automatique - Onglet Paramètres

Comparaison des modèles

1. Cliquez sur le bouton Exécuter.

Le nugget de modèle est créé et placé sur l'espace de travail et dans la palette Modèles en haut à droite de la fenêtre. Vous pouvez parcourir le nugget ou l'enregistrer ou le déployer de plusieurs façons.

Ouvrez le nugget de modèle ; il répertorie les détails concernant chacun des modèles créés au cours de l'exécution. (En situation réelle, lorsque des centaines de modèles sont estimés à partir d'un grand nombre de données, cette opération peut prendre plusieurs heures.) Voir figure 42, à la page 51.

Si vous souhaitez analyser plus en détail l'un des modèles individuels, vous pouvez double-cliquer sur un nugget de modèles dans la colonne **Modèle** pour la faire défiler et parcourir les résultats du modèle individuel ; à partir de là, vous pouvez générer des noeuds de modélisation, des nuggets de modèle ou des graphiques d'évaluation.

😨 taxable_value 🛛 🛛 🕅							
*	File 🖏 Generate 🖋 View 🕞 Preview 🚳						
Model Grap	oh Summary Settings	Annotations					
Sort by:	Correlation 🔹 🔍 A	scending 🔘 Descending		ete Unused Models	View: Training set	-	
Use?	Graph	Model	Build Time (mins)	Correlation 🗸	No. Fields Used	Relative Error	
	Line and the second second	Generalized Linear 1	<1	0.915	7	0.162	
	and the first	Regression 1	<1	0.9	5	0.19	
	CHAID Tree 1 <1 0.892 5 0.204						
ОКСа	ncel					Apply Reset	

Figure 48. Résultats de la numérisation automatique

Par défaut, les modèles sont classés en fonction de la corrélation, cette mesure ayant été sélectionnée dans le noeud Numérisation automatique. Pour faciliter le classement, la valeur absolue de la corrélation est utilisée, avec les valeurs les plus proches de 1 indiquant une relation très forte. Le modèle linéaire généralisé est classé comme étant le meilleur en fonction de cette mesure, mais plusieurs autres sont presque aussi précis. Ce modèle linéaire généralisé a également l'erreur relative la plus basse.

Vous pouvez effectuer le tri sur une autre colonne en cliquant sur l'en-tête de cette colonne ou vous pouvez choisir la mesure désirée dans la liste **Trier par** de la barre d'outils.

Chaque graphique présente un tracé de valeurs observées par rapport aux valeurs prédites pour le modèle et fournit ainsi une indication visuelle rapide de leurs corrélation. Pour un modèle performant, les points doivent être regroupés le long de la diagonale, ce qui est vrai pour tous les modèles de cet exemple.

Dans la colonne **Graphique**, vous pouvez double-cliquer sur une miniature pour générer un graphique en grandeur nature.

En fonction de ses résultats, vous pouvez décider d'utiliser les trois modèles les plus précis. En combinant les prévisions à partir de plusieurs modèles, il est possible d'éviter les limitations dans les modèles individuels. Ce qui entraîne une plus grande exactitude globale.

Dans la colonne Utiliser ?, vérifiez que les trois modèles sont sélectionnés.

Liez un noeud Analyse (palette Sortie) après le nugget de modèle. Cliquez avec le bouton droit de la souris sur le noeud Analyse et sélectionnez **Exécuter** pour exécuter le flux.

La moyenne du score généré par le modèle d'ensemble est ajoutée à un champ *\$XR-taxable_value*, avec une corrélation de 0,922, ce qui est supérieur à ceux des trois modèles individuels. Les scores d'ensemble affichent également une faible erreur absolue moyenne et peuvent être plus efficaces que tous les modèles

individuels lorsqu'ils sont appliqués à d'autres jeux de données.

🛾 Analysi	s of [taxable_value]		×
😂 <u>F</u> ile 🛛	🍃 File 📄 Edit 🔛 🕒 📢		
Analysis	Annotations		
8 Collaps	e All 🍖 Expand All		
Results	for output field taxable_value		
🖻 Con	paring \$XR-taxable_value wit	th taxable_value	
	Minimum Error	-156049.854	
	Maximum Error	176856.403	
	Mean Error	0.014	
	Mean Absolute Error	21353.824	
	Standard Deviation	30815.028	
	Linear Correlation	0.922	
	Occurrences	1,138	
			бк

Figure 49. Flux d'échantillons Numérisation automatique

Récapitulatif

Pour résumer, vous avez utilisé le noeud Numérisation automatique pour comparer plusieurs modèles différents, vous avez sélectionné les trois modèles les plus précis et vous les avez ajoutés au flux dans un nugget de modèle Numérisation automatique combiné.

- Concernant l'exactitude globale, les modèles linéaires généralisés, de Régression et CHAID sont plus performants avec les données d'apprentissage.
- Le modèle d'ensemble a été plus performant que deux des trois modèles individuels et peut être aussi efficace lorsqu'il est appliqué à d'autres jeux de données. Si votre objectif est d'automatiser autant que possible le processus, cette approche vous permet d'obtenir un modèle fiable dans la plupart des circonstances sans avoir à creuser trop dans les spécificités des modèles.

Chapitre 6. Préparation automatique de données (ADP)

La préparation des données pour l'analyse est une des étapes les plus importantes des projets et généralement, l'une de celles qui prend le plus de temps. Le noeud de préparation automatique de données (ADP) gère cette tâche pour vous en analysant vos données et en identifiant des corrections, en filtrant des champs problématiques et peu susceptibles d'être utiles et en créant de nouveaux attributs le cas échéant, et enfin en améliorant la performance au moyen de techniques de filtrage intelligentes. Vous pouvez utiliser le noeud de manière totalement automatisée, en laissant le noeud choisir et appliquer les corrections, ou vous pouvez prévisualiser les modifications avant qu'elles ne soient effectuées et les accepter ou les refuser au choix.

Le noeud ADP vous permet de préparer rapidement et facilement les données pour le Data Mining sans connaissance préalable des concepts statistiques impliqués. Si vous exécutez le noeud avec les paramètres par défaut, les modèles auront tendance à être créé et à réaliser des évaluations plus rapidement.

Cet exemple utilise le flux nommé *ADP_basic_demo.str*, qui se rapporte à un fichier de données nommé *telco.sav* pour expliquer l'exactitude accrue dont vous pouvez bénéficier en utilisant les paramètres par défaut du noeud ADP lors de la création de modèles. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *ADP_basic_demo.str* se trouve dans le répertoire des *flux*.

Création du flux

1. Pour créer le flux, ajoutez un noeud source Statistics qui pointe sur *telco.sav*, dans le répertoire *Demos* du dossier d'installation d'IBM SPSS Modeler.



Figure 50. Création du flux

2. Attachez un noeud type au noeud source, définissez le niveau de mesure du champ *attrition* sur **Indicateur** et le rôle sur **Cible**. Le rôle de tous les autres champs doit être défini sur **Entrée**.

R	Preview				0-1
Types Forme	Appotations				
\ + 000	Read Va	ilues Clear	Values	Clear All Va	lues
Field -	Measurement	Values	Missing	Check	Role
		0.0,1.0		None	
logiong	Continuous	[-0.10536		None	a input
	Continuous	[1.74313		None	Input
logequi	Continuous	[1.01160		None	> Input
logwire	Continuous	[2 70136		None	Input
		[2:19722		None	N Input
custcat	Nominal	1.0.2.0.3		None	
	U Eleve	10/00		None	O Target

Figure 51. Sélection de la cible

- 3. Reliez un noeud Logistique au noeud type.
- 4. Dans le noeud Logistique, cliquez sur l'onglet Modèle et sélectionnez la procédure **Binomial**. Dans le champ *Nom de modèle*, sélectionnez **Personnalisé** et saisissez Pas de ADP attrition.

😡 No ADP - churn		×				
Madel peme: Auto O (ustom	No ADR - churp				
	Justom					
Use partitioned data						
Build model for each split						
Procedure: 🔘 Multinomia	d	Binomial				
Binomial Procedure						
Method: Enter 💌						
Categorical Inputs:						
Field Name	Contrast	Base Category				
		×				
OK Run Cancel		Apply Reset				

Figure 52. Choix des options de modèle

- 5. Reliez un noeud ADP au noeud type. Dans l'onglet Objectifs, conservez les paramètres par défaut afin d'analyser et de préparer vos données en équilibrant la vitesse et l'exactitude.
- 6. En haut de l'onglet Objectifs, cliquez sur **Analyser les données** afin d'analyser et de traiter vos données.

D'autres options du noeud ADP vous permettent de spécifier si vous souhaitez vous concentrer davantage sur l'exactitude, sur la vitesse de traitement ou affiner les nombreuses étapes de traitement de la préparation des données.



Figure 53. Objectifs ADP par défaut

Les résultats du traitement des données sont affichés dans l'onglet Analyse. Le **Récapitulatif de traitement des champs** montre que parmi les 41 éléments de données que propose le noeud ADP, 19 ont été transformés afin d'améliorer le traitement et 3 ont été abandonnés car ils ne sont pas utilisés.

😵 Auto Data Prep							
Cenerate of View Preview Analyze Data Clear Analysis							
Objectives Fields Settings Analysis Annotat	tions						
Field Processing Summary		Predictors Recommended for Use in Analysis Predictive Power					
Fields		N	Target: churn				
<u>Larget</u>		1					
Predictors		41	tenure transformed Equipment				
	Total	38	Internet				
Predictors recommended for use in analysis	Original fields (untransformed)	19	transformed Electronic				
	Transformations of original fields	19	Calling card				
	Derived from dates and times	0	education agetransformed				
	Constructed	0	transformed Customer_				
Predictors not used		3	0.0 0.2 0.4				
			Least Important	Most Important			
View: Field Processing Summary T Reset							
OK Cancel				Apply Reset			

Figure 54. Récapitulatif du traitement des données

- 7. Reliez un noeud Logistique au noeud ADP.
- 8. Dans le noeud Logistique, cliquez sur l'onglet Modèle et sélectionnez la procédure **Binomial**. Dans le champ *Nom de modélisation*, sélectionnez **Personnalisé** et saisissez Après ADP attrition.

🙀 After ADP - churn		X			
Fields Model Expert	Analyze Annotations				
Model name: 🛛 Auto 🧕	Custom	After ADP - churn			
Use partitioned data					
👿 Build model for each sp	lit				
Procedure: 🔘 Multinom	Binomial				
Binomial Procedure					
Method: Enter					
Categorical Inputs:					
Field Name	Contrast	Base Category			
	- 25	×			
📝 Include constant in equ	ation				
OK 🕨 Run Canc	el	Apply Reset			

Figure 55. Choix des options de modèle

Comparaison de l'exactitude des modèles

1. Exécutez les deux noeuds Logistique pour créer les nuggets de modèle, qui sont ajoutés au flux et à la palette Modèles dans l'angle supérieur droit.



Figure 56. Relier les nuggets de modèle

2. Reliez les noeuds Analyse aux nuggets de modèle et exécutez des noeuds Analyse avec leurs paramètres par défaut.



Figure 57. Relier les noeuds d'analyse

L'analyse du modèle non dérivé ADP montre que la seule exécution des données dans le noeud Régression logistique avec ces paramètres par défaut fournit un modèle de faible exactitude - seulement 10,6 %.

No ADP	- LogReg						
😂 <u>F</u> ile	🖻 Edit 🛛 🔞		1	@ ×			
Analysis	Annotations	536° - 6					
Collapse All 🌾 Expand All							
-Results	■-Results for output field churn						
🖻 Cor	nparing \$L-churr	n with chur	'n				
	Correct	106	10.6%				
	Wrong	894	89.4%				
	Total	1,000	energine ne ne sé				
				OK			

Figure 58. Résultat d'un modèle non dérivé de l'ADP

L'analyse du modèle dérivé de l'ADP montre que dans le cadre de l'exécution des données avec les paramètres ADP par défaut, vous avez construit un modèle beaucoup plus précis, exact à 78.8%.
Image: Second and Second	×
Analysis Annotations & Collapse All P Expand All ••• Results for output field churn •••• Comparing \$L-churn with churn •••• Correct 788 78.8% Wrong 212 21.2% Total 1,000	
Collapse All Expand All Correct 788 78.8% Wrong 212 21.2% Total 1,000	
Results for output field churn Comparing \$L-churn with churn Vorrect 788 78.8% Wrong 212 21.2% Total 1,000	
Comparing \$L-churn with churn Correct 788 78.8% Wrong 212 21.2% Total 1,000	
Correct 788 78.8% Wrong 212 21.2% Total 1,000 1	
Wrong 212 21.2% Total 1,000	
Total 1,000	
<u></u>	

Figure 59. Résultat d'un modèle dérivé de l'ADP

Dans le récapitulatif, en exécutant uniquement le noeud ADP pour affiner le traitement de vos données, vous avez été en mesure de construire un modèle plus précis avec peu de manipulation directe des données.

Bien sûr, si votre objectif est de prouver ou non la validité d'une certaine théorie, ou si vous souhaitez construire des modèles spécifiques, il peut être préférable d'utiliser directement les paramètres de modèle. Cependant, pour les personnes disposant de peu de temps ou si vous avez de grandes quantités de données à préparer, le noeud ADP peut représenter un avantage.

Des explications sur les fondements mathématiques des méthodes de modélisation utilisées dans IBM SPSS Modeler sont présentées dans le *guide des algorithmes d'IBM SPSS Modeler*, disponible dans le répertoire *Documentation* du disque d'installation.

Sachez également que les résultats de cet exemple sont établis uniquement d'après les données d'apprentissage. Pour évaluer comment les modèles peuvent se généraliser à d'autres données dans le monde réel, vous devez utiliser un noeud Partitionner destiné à contenir un sous-ensemble d'enregistrements à des fins de test et de validation.

Chapitre 7. Préparation des données pour l'analyse (Audit données)

Le noeud Audit données fournit un premier aperçu complet des données importées dans IBM SPSS Modeler. Souvent utilisé lors de l'exploration initiale des données, le rapport d'audit des données affiche des statistiques récapitulatives, ainsi que les histogrammes et les graphiques de distribution pour chaque champ de données. Il vous permet en outre d'indiquer comment traiter les valeurs manquantes, les valeurs éloignées et les valeurs extrêmes.

Cet exemple utilise le flux *telco_dataaudit.str*, qui fait référence au fichier de données *telco.sav*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *telco_dataaudit.str* se trouve dans le répertoire des *flux*.

Création du flux

1. Pour créer le flux, ajoutez un noeud source Statistics qui pointe sur *telco.sav*, dans le répertoire *Demos* du dossier d'installation d'IBM SPSS Modeler.



Figure 60. Création du flux

Ajoutez un noeud type pour définir des champs, puis désignez *attrition* comme champ cible (Rôle = Cible). Le rôle doit avoir la valeur Entrée pour tous les autres champs pour que cette cible soit la seule cible.

•>	Preview			1	0-1
Types Forma	t Annotations				
	Read Va	liues Clear	values	Clear All Va	alues
Field -	Measurement	Values	Missing	Check	Role
	Continuouo	1/0		None	Input
Togiong	Continuous	[-0.10536		None	a input
		[1.74919		None	= input
rogequi	Continuous	[2.73436		INORE	input
📸 logcard	Continuous	[1.01160	_	None	> Input
nogwire 👷	Continuous	[2.70136		None	🔪 Input
ninc 🖉	🖉 Continuous	[2.19722		None	🔪 Input
💭 custcat	🎳 Nominal	1,2,3,4		None	🔪 Input
	🙎 Elag	1/0		None	O Target

Figure 61. Définition de la cible

3. Vérifiez que les niveaux de mesure de champ sont correctement définis. Par exemple, la plupart des champs dont les valeurs sont 0 et 1 peuvent être considérés comme des champs indicateurs. Cependant, certains champs, tels que celui indiquant le genre, doivent être considérés comme des champs nominaux à deux valeurs.

3 DE	review				0.
<u> </u>					
Types Format	Annotations				
% -	💌 🚺 🕨 Read Va	alues Clear	Values	Clear All Va	lues
Field -	Measurement	Values	Missing	Check	Role
🏓 ed	📲 Ordinal	1,2,3,4,5	_	None	🔪 Input
employ	🔗 Continuous	[0,47]		None	🔪 Input
🔉 retire	💑 Nominal	0.0,1.0		None	🔪 Input
> gender	💑 Nominal	0,1		None	🔪 Input
🔉 reside	🚮 Ordinal	1,2,3,4,5,		None	🔪 Input
🔉 tollfree	🖁 Flag	1/0		None	🔪 Input
equip	🎖 Flag	1/0		None	🔪 Input
Callcard	🎖 Flag	1/0		None	🔪 Input
> wireless	🙎 Flad	1/0		None	🔪 Innut
View current	fields 🔘 View upu	read field cattin	10		
y view current		iseu neia settini	38		

Figure 62. Définition des niveaux de mesure

Astuce : Pour modifier les propriétés de plusieurs champs contenant des valeurs similaires (telles que 0/1), cliquez sur l'en-tête de colonne *Valeurs* afin de trier les champs en fonction de cette colonne. Utilisez la touche Maj pour sélectionner tous les champs à modifier. Cliquez ensuite sur la sélection avec le bouton droit de la souris pour modifier le niveau de mesure ou les autres attributs de tous les champs sélectionnés.

4. Connectez un noeud Audit données au flux. Dans l'onglet Paramètres, conservez les paramètres par défaut pour que tous les champs soient inclus dans le rapport. Etant donné que *attrition* est le seul champ cible défini dans le noeud type, ce champ est automatiquement utilisé comme champ de superposition.

😡 42 Fie	lds						X
						0	
Settings	Quality	Output	Annotations				
🔘 Default			O U:	se custo	om fields		
Fields:							×
Overlay:							-
Display Grap	hs	∎ В	asic statistics		🗖 Adva	nced statisti	CS
🗾 Calcula	te mediar	n and mod	e (may slow p	erforma	nce on larg	je datasets)	
ОК	Run	Cancel				Appl	y <u>R</u> eset

Figure 63. Noeud Audit des données - Onglet Paramètres

Dans l'onglet Qualité, conservez les paramètres par défaut de détection des valeurs manquantes, éloignées et extrêmes, puis cliquez sur **Exécuter**.

42 Fie	elds				0 -
ettings	Quality	Output	Annotations		
Calculate Calculate Cour Brea Dutliers & Detect	alues at of recol kdown co Extreme tion Metho andard da	rds with v punts of re Values – od: eviation fro	alid values scords with inv	alid values	
Out	liers:	3.0 🗲	Extremes:	5.0 ≑	
() Inte	erquartile	ranges fr	om upper/lowe	r quartiles	
Out	liers:	1.5	Extremes:	3.0 🌲	
Note: \$	Selecting	Interquart	ile range may s	low performanc	e on large datasets
	Burn	Connall			

Figure 64. Noeud Audit des données - Onglet Qualité

Navigation dans les statistiques et les graphiques

Le navigateur Audit des données est affiché avec des graphiques en miniature et des statistiques descriptives pour chaque champ.

🔍 Data Au	dit of [42 fields]								
🐞 <u>F</u> ile 🔋	🛓 Edit 🛛 🕙 Generate								?
Audit Qual	lity Annotations								
Field -	Graph	Measurement	Min	Max	Mean	Std. Dev	Skewness	Unique	Valid
🔷 region		💑 Nominal	1	3			-	3	1000
🔷 tenure		🛷 Continuous	1	72	35.526	21.360	0.112		1000
🔷 age		🔗 Continuous	18	77	41.684	12.559	0.357		1000
🔷 marital		🎖 Flag	o	1			-	2	1000
🔷 address		🔗 Continuous	o	55	11.551	10.087	1.106		1000
🋞 income	l	🔗 Continuous	9.000	1668.000	77.535	107.044	6.643	8	1000
1 Indicates a r	nuttimode result 2 Indi	cates a sampled result							ок

Figure 65. Navigateur Audit des données

A l'aide de la barre d'outils, affichez les libellés de champ et de valeur, et basculez l'alignement des graphiques de l'horizontale à la verticale (champs catégoriels uniquement).

1. La barre d'outils ou le menu Editer vous permet en outre de choisir les statistiques à afficher.

		1
	Statistic	
-	Min	
-	Max	
	Sum	
	Range	
-	Mean	
	Mean Std. Err.	
-	Standard deviation	
10	Variance	
-	Skewness	
	Skewness Std. Err.	
	Kurtosis	
	Kurtosis Std. Err.	
-	Unique	
-	Valid	
	OK Help	_

Figure 66. Affichage des statistiques

Double-cliquez sur un graphique en miniature dans le rapport d'audit pour afficher ce graphique en taille réelle. Etant donné que *Flux* est le seul champ cible du flux, il est automatiquement utilisé comme champ de superposition. Vous pouvez basculer l'affichage des libellés de champ et de valeur à l'aide de la barre d'outils de la fenêtre Graphiques ou cliquer sur le bouton Mode d'édition pour personnaliser le graphique.



Figure 67. Histogramme de durée d'affectation

Vous pouvez également sélectionner une ou plusieurs miniatures et générer un noeud graphique pour chacune d'elles. Les noeuds générés sont placés dans le canevas de flux. Vous pouvez les ajouter au flux pour recréer le graphique concerné.

s SuperNode eme SuperNode s Filter Node s Select Node de	Min 1 1 1	<u>Мах</u> 3 72	Mean	Std. Dev 21.360	Skewness	Unique 3 	✓alid 1000
eme SuperNode s <u>Filter Node</u> s Select Node de	Min 1	Max 3 72	Mean 35.526	Std. Dev 21.360	Skewness	Unique 3 	Valid 1000
s Eilter Node s Select Node de	1	3	35.526	21.360	0.112	3	1000
	1	72	35.526	21.360	0.112		1000
2	18						
<i>.</i> 0		77	41.684	12.559	0.357		1000
🎖 Flag	0	1	77			2	1000
Continuous	O	55	11.551	10.087	1.106	6. 	1000
🖉 Continuous	9.000	1668.000	77.535	107.044	6.643	6 <u>11</u>	1000
	Continuous Continuous stes a sampled result		Continuous Continuous Soud Soud	Continuous O Continuous O S O S S O S S O S S O S S O S	Continuous 0 55 11.551 10.087 Continuous 9.000 1668.000 77.535 107.044 ates a sampled result V	Continuous 0 55 11.551 10.087 1.106 Continuous 9.000 1668.000 77.535 107.044 6.643 ates a sampled result	Continuous 0 55 11.551 10.087 1.106 Continuous 9.000 1668.000 77.535 107.044 6.643 ates a sampled result

Figure 68. Génération d'un noeud graphique

Traitement des valeurs éloignées et manquantes

L'onglet Qualité du rapport d'audit contient des informations sur les valeurs éloignées, extrêmes et manquantes.

udit Quality Ann	otationa					
	lotations					
mplete fields (%):	90.476190 Comp	lete records (%):	13.1			
Field 🖘 🛛 🕅	1easurement	Outliers	Extremes	Action	Impute Missing	Method
🕨 region 🛛 🕹	Nominal	222	-		Never	Fixed
🕻 tenure	Continuous	0	0	None	Never	Fixed
bage 🔗	Continuous	0	C	None	Never	Fixed
🕻 marital	Flag		-		Never	Fixed
🔪 address 👘 🔗	Continuous	12	0) None	Never	Fixed
🕻 income 🛛 🔗	Continuous	9	e	None	Never	Fixed
ed 📶	Ordinal		-		Never	Fixed
employ 🧳	Continuous	8	C) None	Never	Fixed
🕻 retire 🛛 🍰	Nominal	8 ()	-		Never	Fixed
👌 gender 🛛 🕹	Nominal	822	2.		Never	Fixed
>reside 🚮	Ordinal		-		Never	Fixed
🕻 tollfree	Flag	842	2		Never	Fixed
🕻 equip	Flag	822	2.	21	Never	Fixed
Callcard 🎖	Flag	822	2.		Never	Fixed
🕻 wireless	Flag	228	2.		Never	Fixed
🕻 longmon 🛛 🗳	Continuous	18	4	None	Never	Fixed
tollmon 🧳	Continuous	9	1	None	Never	Fixed
equipmon 🔗	Continuous	2	C) None	Never	Fixed
cardmon 🔗	Continuous	11	3	None	Never	Fixed

Figure 69. Navigateur Audit de données - Onglet Qualité

Vous pouvez également définir des méthodes de gestion des valeurs et générer des super noeuds qui appliquent automatiquement les transformations. Par exemple, vous pouvez sélectionner un ou plusieurs champs et choisir d'attribuer ou de remplacer les valeurs manquantes de ces champs à l'aide de diverses méthodes, dont l'algorithme C&RT.

Eile 🌛 Eo	lit 🕙 <u>G</u> enerate						ð
udit Quality	Annotations						
omplete fields (%): 90.47619(Com	plete records (%):	13.1				
Field	Measurement	Outliers	Extremes	Action	Impute Missing	Method	0
region	Nominal				Never	Fixed	-
tenure	Continuous	0		0 None	Never	Fixed	-
age	Continuous	0		0 None	Never	Fixed	-
marital	& Flag	120			Blank & Null Values	Fixed 💌	
address	Continuous	12		0 None	Never	Fixed	
income	Continuous	9		6 None	Never	Random	
ed	Ordinal	1. A A A A A A A A A A A A A A A A A A A			Never	Evpression	
employ	Continuous	8		0 None	Never	Algorithm	-
retire	Nominal	844			Never	Specify R	
gender	& Nominal	822		22 220	Never	Fixed	
reside	Ordinal	822			Never	Fixed	
tollfree	🖁 Flag	842			Never	Fixed	
equip	🔓 Flag	822			Never	Fixed	
callcard	🖁 Flag	842			Never	Fixed	
wireless	🖁 Flag	842			Never	Fixed	
longmon	Continuous	18		4 None	Never	Fixed	
tollmon	Continuous	9		1 None	Never	Fixed	
equipmon	Continuous	2		0 None	Never	Fixed	
cardmon	Continuous	11		3 None	Never	Fixed	
					According		

Figure 70. Choix d'une méthode d'attribution

Après avoir spécifié une méthode d'attribution pour un ou plusieurs champs, pour générer un super noeud Valeurs Manquantes, dans les menus choisissez :

Générer >	Super	noeud	des	valeurs	manquantes
-----------	-------	-------	-----	---------	------------

🔰 File 🛛 📄 Eo	dit 👋 <u>G</u> enerate					0	Ĩ
Audit Quality	Missing Values	SuperNode					
- Cicili	Outlier & Extrem	ie SuperNode					
Complete fields (%): Missing Values	Filter blade	1				
	wissing values	Eliter Node					
Field -	Missing Values	Select Node	xtremes	Action	Impute Missing	Method	
> region	Reclassify Node	2			Never	Fixed	-
tenure	Ø Disaise Meda		() None	Never	Fixed	
age	Binning Node	0 None			Never	Fixed	
🕻 marital	Derive Node		-		Blank & Null Values	Fixed	-
address	Croph Output		() None	Never	Fixed	
income	Sraph Output		6	3 None	Never	Fixed	
> ed	Graph Node		-		Never	Fixed	
> employ	🞸 Continuous	8	- () None	Never	Fixed	
retire	💑 Nominal		-		Never	Fixed	
> gender	💑 Nominal	842	-		Never	Fixed	
> reside	I Ordinal	842 ()	-		Never	Fixed	
> tollfree	🥉 Flag	822	-	-	Never	Fixed	
👌 equip	🍯 Flag	842 () 	-		Never	Fixed	
> callcard	🥉 Flag	822			Never	Fixed	
> wireless	🥉 Flag	8-2			Never	Fixed	
longmon	Continuous	18		4 None	Never	Fixed	
tollmon	Continuous	9	1	1 None	Never	Fixed	
equipmon	Continuous	2	() None	Never	Fixed	
🖗 cardmon	🞸 Continuous	11		3 None	Never	Fixed	
and the state of t				and the second se			

Figure 71. Génération du super noeud

Le super noeud généré est ajouté au canevas de flux, où vous pouvez le connecter au flux pour appliquer les transformations.



Figure 72. Flux avec super noeud Valeurs manquantes

Le super noeud contient en réalité plusieurs noeuds qui exécutent les transformations requises. Pour comprendre son fonctionnement, modifiez le super noeud et cliquez sur **Zoom avant**.



Figure 73. Zoom avant sur le super noeud

Chaque champ auquel une valeur est attribuée à l'aide de la méthode algorithmique, par exemple, est associé à un modèle C&RT distinct et à un noeud Remplacer qui remplace les valeurs non renseignées et les valeurs nulles par la valeur prédite par le modèle. Vous pouvez ajouter, modifier ou supprimer des noeuds précis dans le super noeud pour personnaliser encore davantage son comportement.

Vous pouvez également générer un noeud Sélectionner ou Filtrer pour supprimer les champs ou les enregistrements où des valeurs manquent. Par exemple, vous pouvez filtrer les champs dont le pourcentage de qualité est inférieur au seuil défini.



Figure 74. Génération d'un noeud Filtrer

Les valeurs éloignées et extrêmes peuvent être gérées de manière similaire. Indiquez l'action à appliquer à chaque champ (Forcer, Supprimer ou Rendre nul) et générez un super noeud pour appliquer les transformations.

File 📴 Edit	O Generate					0	
utit Quality	Missing Values Su	perNode					
acate and and a second se	Outlier & Extreme :	SuperNode					
mplete fields (%): Missing Values <u>F</u> il	er Node	\$ 1				
Field	Missing Values Se	lect Node	xtremes	Action	Impute Missing	Method	
region	Reclassify Node		-		Never	Fixed	
tenure .	Disasing March		0	None	Never	Fixed	
age .	Einning Node		0	None	Never	Fixed	
marital	Derive Node				Never	Fixed	
address	Overski Overset		0	Coerce	Blank & Null Val	Fixed	
income .	Graph Output		6	None	Never	Fixed	
ed	Graph Node				Never	Fixed	
employ .	🖉 Continuous	8	0	None	Never	Fixed	
retire	💑 Nominal				Never	Fixed	
gender	💑 Nominal	822			Never	Fixed	
reside	📶 Ordinal	8442		120	Never	Fixed	
tollfree	🎖 Flag	87 -1 2			Never	Fixed	
equip	🎖 Flag	8922		<u>1</u> 1	Never	Fixed	
callcard	🎖 Flag	8442			Never	Fixed	
wireless	🎖 Flag	85 -2 2	<u></u>	<u></u>	Never	Fixed	
longmon .	🔗 Continuous	18	4	None	Never	Fixed	
tollmon .	🔗 Continuous	9	1	None	Never	Fixed	
equipmon	🔗 Continuous	2	0	None	Never	Fixed	
cardmon .	🔗 Continuous	11	3	None	Never	Fixed	
wiremon .	🔗 Continuous	8	1	None	Never	Fixed	
longten .	🔗 Continuous	20	4	None	Never	Fixed	
tollten .	🔗 Continuous	18	2	None	Never	Fixed	
equipten .	🔗 Continuous	16	3	None	Never	Fixed	
cardten	🔗 Continuous	11	6	None	Never	Fixed	
wireten .	🖉 Continuous	22	3	None	Never	Fixed	
multline	🎖 Flag				Never	Fixed	
				and the second se			

Figure 75. Génération d'un noeud Filtrer

Une fois l'audit terminé et les noeuds générés ajoutés au flux, vous pouvez poursuivre l'analyse. Vous pouvez également effectuer une analyse plus poussée des données grâce à la méthode Détection des anomalies ou Sélection de fonction, ou à d'autres méthodes.



Figure 76. Flux avec super noeud Valeurs manquantes

Chapitre 8. Traitements par médicaments (Graphiques exploratoires/C5.0)

Pour cette section, imaginez que vous êtes un chercheur et que vous souhaitez compiler des données pour une étude médicale. Vous avez rassemblé des données sur un ensemble de patients, souffrant tous de la même maladie. Lors du traitement, chaque patient a réagi à l'un des cinq médicaments. Votre travail consiste à utiliser l'exploration de données pour savoir quel médicament pourrait convenir à un futur patient atteint de la même maladie.

Cet exemple utilise le flux intitulé *druglearn.str*, qui référence le fichier de données *DRUG1n*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *druglearn.str* se trouve dans le répertoire des *flux*.

Champ de données	Description
Age	(Nombre)
Sexe	M ou F
ТА	Tension artérielle : HIGH (ELEVEE), NORMAL ou LOW (BASSE)
Cholestérol	Taux de cholestérol dans le sang : NORMAL ou HIGH
Na	Concentration de sodium dans le sang
K	Concentration de potassium dans le sang
Médicament	Médicament prescrit auquel le patient a réagi

Les champs de données utilisés dans la démo sont :

Lecture de données texte



Figure 77. Ajout d'un noeud Délimité

Vous pouvez lire des données texte délimitées à l'aide d'un **noeud Délimité**. Vous pouvez ajouter un noeud Délimité depuis les palettes en cliquant sur l'onglet **Sources** pour rechercher le noeud ou utiliser l'onglet **Favoris** qui contient ce noeud par défaut. Ensuite, double-cliquez sur le noeud que vous venez de placer pour ouvrir la boîte de dialogue correspondante.

Pour sélectionner le répertoire dans lequel IBM SPSS Modeler est installé sur votre système, cliquez sur le bouton représentant des points de suspension (...), à droite de la zone Fichier. Ouvrez le répertoire *Demos*, puis sélectionnez le fichier *DRUG1n*.

Vérifiez que vous avez sélectionné **Lire noms des champs à partir du fichier** et notez les valeurs et champs qui ont été chargés dans la boîte de dialogue.

😡 Var. File	
\$CLEO_DEMOS(DROOTH	
File Data Filter Types Annotations	
File: \$CLEO_DEMOS\DRUG1n	
Age Sey PD Cholesterol Ne K Du	
23,F,HIGH,HIGH,0.792535,0.0312	258, drugY
47,M,LOW,HIGH,0.739309,0.05640	68, drugC
1,7,17,10,00,7,10,7,00,7,10,7,00,00,7,10,00,7,10,00,00,00,00,00,00,00,00,00,00,00,00,	
Read field names from file	Specify number of fields
Skip header characters: 0 🗧	EOL comment characters:
Strip lead and trail spaces: 🛛 🔘 None 🔘 L	eft 🔘 Right 🔘 Both
Invalid characters: 💿 Discard 🔘	Replace with
Encoding: Stream default 🔻	Decimal symbol: Stream default 🔻
Delimiters	Lines to scan for type: 50 🗧
🔄 Space 🛛 🧹 Comma 📃 Tab	Automatically recognize dates and times
Vewline 🗾 Other	-Quotes
Non-printing characters	Single quotes: Discard
Allow multiple blank delimiters	Double quotes: Discard
OK Cancel	Apply

Figure 78. Boîte de dialogue Délimité

🚰 DRUG1 n			
CLEO_DEMOS/DRUG1	sh I n tions		0
Field	Override	Storage	Input Format
View current fields View un OK Cancel	nused field settin	gs	Apply Reset

Figure 79. Modification du type de stockage pour un champ

😡 DRUG1 n					
	view 😰 Refresh				
\$CLEO_	DEMOS/DRUG1n				
File Data Filter	Types Annotations				
~	Read Values	Clear Values	Clear All	Values	
Field -	Measurement	Values	Missing	Check	Role
🚫 Age	🔗 Continuous	[15,74]		None	🔪 Input
A Sex	🎖 Flag	M/F		None	🔪 Input
A BP	💑 Nominal	HIGH,LOW,		None	🔪 Input
A Cholesterol	🎖 Flag	NORMAL/HI	Off 🛛 💌	None	🔪 Input
🛞 Na	🔗 Continuous	[0.500169,0	On (*)	None	🔪 Input
🛞 К 💦 👘	🔗 Continuous	[0.020022,0	Off K	None	🔪 Input
A Drug	💑 Nominal	drugA,drug	Specify	None	🔪 Input
View current fie	elds 🔘 View unused fie	eld settings			
OK Cancel					Apply Reset

Figure 80. Sélection des options de valeur dans l'onglet Types

Cliquez sur l'onglet **Données** pour ignorer et modifier le **Stockage** d'un champ. Veuillez noter que le stockage est différent des **Mesures**, c'est-à-dire, le niveau de mesure (ou le type d'utilisation) du champ des données. L'onglet **Types** vous permet d'obtenir des informations supplémentaires sur le type de

champs de vos données. Vous pouvez également choisir **Lire les valeurs** pour afficher les valeurs réelles de chaque champ en fonction des sélections effectuées dans la colonne *Valeurs*. Ce processus est appelé **instanciation**.

Ajout d'une table

Maintenant que le fichier de données est chargé, vous pouvez examiner les valeurs de certains enregistrements. Pour ce faire, vous pouvez, par exemple, créer un flux incluant un noeud Table. Pour placer un noeud Table dans le flux, double-cliquez sur l'icône de la palette ou faites-la glisser vers l'espace de travail.



Figure 81. Noeud Table relié à la source de données

<u>File E</u> dit Insert <u>V</u>	iew <u>T</u> ools <u>S</u> uperN	lode <u>Wi</u> n	wob	Help	i.					
				5	~				* *	*
		12						R	un the current stream	
		Table	(7 f	ields	, 200 re	cords) #2				
		in an			#D 0			A 36		
		File	3	Edit		rate 📖				
	_	Table ,	Annota	tions						
B			Age	Sex	BP	Cholesterol	Na	ĸ	Drug	
		1	23	F	HIGH	HIGH	0.793	0.031	drugY	-
		2	47	М	LOW	HIGH	0.739	0.056	drugC	
DRUG1n	Table	3	47	М	LOW	HIGH	0.697	0.069	drugC	
		4	28	F	NORMAL	HIGH	0.564	0.072	drugX	
		5	61	F	LOW	HIGH	0.559	0.031	drugY	
		6	22	F	NORMAL	HIGH	0.677	0.079	drugX	
		7	49	F	NORMAL	HIGH	0.790	0.049	drugY	
		8	41	М	LOW	HIGH	0.767	0.069	drugC	
		9	60	М	NORMAL	HIGH	0.777	0.051	drugY	
		10	43	М	LOW	NORMAL	0.526	0.027	drugY	
		11	47	F	LOW	HIGH	0.896	0.076	drugC	
		12	34	F	HIGH	NORMAL	0.668	0.035	drugY	
		13	43	М	LOW	HIGH	0.627	0.041	drugY	
		14	74	F	LOW	HIGH	0.793	0.038	drugY	
		15	50	F	NORMAL	HIGH	0.828	0.065	drugX	
		16	16	F	HIGH	NORMAL	0.834	0.054	drugY	
		17	69	М	LOW	NORMAL	0.849	0.074	drugX	
		18	43	М	HIGH	HIGH	0.656	0.047	drugA	
		19	23	М	LOW	HIGH	0.559	0.077	drugC	
		20	32	F	HIGH	NORMAL	0.643	0.025	drugY	*
										ОК

Figure 82. Exécution d'un flux à partir de la barre d'outils

lorsque vous double-cliquez sur un noeud de la palette, il est automatiquement relié au noeud sélectionné dans l'espace de travail. Si les noeuds ne sont pas reliés, vous pouvez également utiliser le bouton central de votre souris pour relier le noeud source au noeud Table. Pour simuler l'action du bouton central de la souris, maintenez la touche Alt enfoncée tout en déplaçant la souris. Pour afficher la table, cliquez dans la barre d'outils sur le bouton représentant une flèche verte afin d'exécuter le flux, ou cliquez avec le bouton droit de la souris sur le noeud Table et sélectionnez **Exécuter**.

Création d'un graphique de distribution

Lors de l'exploration de données, il est souvent utile d'explorer les données en créant des récapitulatifs visuels. IBM SPSS Modeler propose plusieurs types de graphiques en fonction du genre de données à récapituler. Par exemple, pour connaître la proportion des patients ayant réagi à chaque médicament, utilisez un noeud distribution.

Ajoutez un noeud distribution au flux, connectez-le au noeud source, puis cliquez deux fois dessus pour éditer les options d'affichage.

Sélectionnez *Médicament* comme champ cible dont vous souhaitez afficher la proportion. Ensuite, cliquez sur le bouton **Exécuter** de la boîte de dialogue.

😡 Dru	ıg			$\overline{\mathbf{X}}$
	Field: Drug			0
Plot	Appearance Output	Annotations		
Plot:	Selected fields	(🕽 All flags (true valu	ies)
Field:	🖁 Drug			
Overl	🖲 🖲 Natural 🔘 Name (O Type	W.	
Color	(none) Sex BP Cholesterol			
Sort:	portional scale	Drug		
ОК	Run Cancel			Apply Reset

Figure 83. Sélection du médicament en tant que champ cible

崖 Distribution of	of Drug #1			
😺 Eile 📄 Edit	🏷 Generate 🛛 💰	⊻iew		0 ×
Table Graph A	nnotations			
Value 🗠	Prop	ortion	%	Count
drugA]		11.5	23
drugB			8.0	16
drugC			8.0	16
drugX			27.0	54
drugY			45.5	91
				ОК

Figure 84. Proportion des réactions à un type de médicament

Le graphique obtenu vous permet de visualiser la forme des données. Il démontre que les patients ont réagi le plus souvent au médicament *Y* et moins souvent aux médicaments *B* et *C*.

🔍 Data Audi	🖳 Data Audit of [7 fields] 📃 🗖 🔀						
🐞 File 🏾 📄	Edit 🛛 🕙 Generate					0 ×	
Audit Quality	Annotations						
Field -	Graph	Measurement	Min	Max	Mean	Std. Dev	
🔷 Age		🔗 Continuous	15	74	44.315	16.544	
A Sex		Sategorical					
A BP		Categorical					
A Cholest		Categorical			80		
🛞 Na	ſŀalaĥmr	🛷 Continuous	0.500	0.896	0.697	0.119	
ж		🔗 Continuous	0.020	0.080	0.050	0.018	
🛕 Drug		Categorical	-	-	-	-	
1	antes fee Managaleman		and the second	California California		•	
¹ Indicates a mu	itimode result ² Indica	ates a sampled result					
						ОК	

Figure 85. Résultats d'un audit de données

Vous pouvez également relier et exécuter un noeud Audit données pour obtenir un aperçu des proportions et des histogrammes de tous les champs simultanément. Le noeud Audit données est disponible dans l'onglet Sortie.

Création d'un nuage de points

A présent, examinons les facteurs susceptibles d'influencer *Médicament*, la variable cible. En tant que chercheur, vous savez que les concentrations de sodium et de potassium dans le sang sont des facteurs importants. Etant donné qu'il s'agit de valeurs numériques, vous pouvez créer un nuage de points pour comparer les valeurs du sodium et du potassium, et utiliser les catégories de médicaments en tant que valeurs de superposition.

Placez un noeud tracé dans l'espace de travail, connectez-le au noeud source, puis cliquez deux fois dessus pour l'éditer.

💽 ? v. ?
X: Na Y: K
Plot Options Appearance Output Annotations
L. X field: Na V field: K J
Panel: Animation: Transparency:
Overlay type: None
O Smoother
Function y =
OK Run Cancel Apply Reset

Figure 86. Création d'un nuage de points

Dans l'onglet Tracé, sélectionnez *Na* comme champ X, *K* comme champ Y et *Médicament* comme champ de superposition. Puis cliquez sur **Exécuter**.

Le tracé montre clairement qu'il existe un seuil au-delà duquel le médicament approprié est toujours le médicament Y et en dessous duquel le médicament approprié n'est jamais le médicament Y. Ce seuil est un rapport : le rapport entre le sodium (*Na*) et le potassium (*K*).



Figure 87. Nuage de points de la proportion de médicaments

Création d'un graphique Relations

La plupart des champs de données étant de type catégoriel, vous pouvez également essayer de tracer un graphique Relations qui réalise le mappage des associations entre les différentes catégories. Connectez un noeud Relations au noeud source dans l'espace de travail. Dans la boîte de dialogue du noeud Relations, sélectionnez *TA* (pression artérielle) et *Médicament*. Puis cliquez sur **Exécuter**.



Figure 88. Graphique Relations médicaments et tension artérielle

Dans le graphique, il semble que le médicament *Y* soit associé aux trois niveaux de pression artérielle. Ceci n'est pas surprenant : vous aviez déjà déterminé dans quel cas du médicament *Y* est approprié. Pour étudier les autres médicaments, vous pouvez masquer le médicament *Y*. Dans le menu **Vue**, choisissez **Mode d'édition**, puis cliquez avec le bouton droit de la souris sur le médicament *Y* et choisissez **Masquer et redessiner**.

Dans le graphique simplifié, le médicament Y et tous ses liens sont masqués Maintenant, vous pouvez voir clairement que seuls les médicaments A et B sont associés à une pression artérielle élevée. Seuls les médicaments C et X sont associés à une pression artérielle faible. Seul le médicament X est associé à une pression artérielle normale. Cependant, vous ne savez toujours pas comment choisir entre le médicament A et le médicament B, ou entre les médicaments C et X, pour un patient donné. C'est dans un cas comme celui-ci que la modélisation peut vous aider.



Figure 89. Graphique Relations avec le médicament Y et ses liens masqués

Calcul d'un nouveau champ

Etant donné que le rapport sodium/potassium semblait indiquer quand utiliser le médicament *Y*, vous pouvez calculer un champ contenant la valeur de ce rapport pour chaque enregistrement. Ce champ peut s'avérer utile par la suite, lors de la création d'un modèle permettant de savoir quand utiliser chacun des cinq médicaments. Pour simplifier la présentation du flux, commencez par effacer tous les noeuds à l'exception du noeud source DRUG1. Reliez un noeud dériver (Onglet Ops sur champs) à DRUG1n, puis cliquez deux fois sur le noeud dériver pour le modifier.

😡 Deriv	e		
	Preview		0
+-5	Derive as: Formula		
Settings	Annotations		
	Mode	: 💿 Single 🔘 Multiple	
Derive fie	eld:		
Na_to_K			
Derive as: Field type Formula:	Formula 💌 8: 🗳 <default> 💌</default>		
NaK			
ОК Са	ancel		Apply Reset

Figure 90. Edition du noeud dériver

Appelez le nouveau champ *Na_sur_K*. Etant donné que vous obtenez le nouveau champ en divisant la valeur du sodium par la valeur du potassium, entrez Na/K dans le champ Formule. Vous pouvez également créer une formule en cliquant sur l'icône située juste à droite du champ. Le Générateur de formules apparaît. Il permet de créer des formules de façon interactive en utilisant des listes de fonctions intégrées, des opérandes, ainsi que des champs et leurs valeurs.

Vous pouvez observer la proportion du nouveau champ en reliant un noeud Histogramme au noeud Calculer. Dans la boîte de dialogue du noeud Histogramme, indiquez que *Na_sur_K* constitue le champ à reporter et *Médicament* le champ de superposition.

?				2
				0-
😬 х:	Na_to_K			
Plot Option	s Appearance	Output	Annotations	
Field: 岁 N	a_to_K			
Coverlay				
Color:	Dr 🚽 Pa	nel:	4	Animation:
OK 🕨 F	tun Cancel			Apply Rese

Figure 91. Edition du noeud Histogramme

Lorsque vous exécutez le flux, vous obtenez le graphique affiché ici. En fonction des éléments affichés, vous pouvez conclure que lorsque la valeur Na_sur_K est égale ou supérieure à 15, le médicament recommandé est le médicament Y.



Figure 92. Histogramme

Création d'un modèle

En exploitant et en manipulant les données, vous avez formulé des hypothèses. Le rapport entre le sodium et le potassium dans le sang semble influer sur le choix du médicament, tout comme la pression artérielle. Mais vous ne pouvez pas encore expliquer totalement tous les liens existant entre ces facteurs. La modélisation vous fournira probablement des réponses. Pour cela, vous pouvez essayer d'ajuster les données à l'aide d'un modèle de création de règle, le modèle C5.0.

Etant donné que vous utilisez un champ déjà calculé, *Na_sur_K*, vous pouvez filtrer les champs d'origine, *Na* et *K*, afin d'éviter qu'ils soient utilisés deux fois dans l'algorithme de modélisation. Vous pouvez utiliser pour cela un noeud Filtrer.

🜍 Filter		$\overline{\mathbf{X}}$
		0.0
Filter Annotations		
7.	Fields:	8 in, 2 filtered, 0 renamed, 6 out
Field -	Filter	Field
Age	\rightarrow	Age
Sex	\rightarrow	Sex
BP	\rightarrow	BP
Cholesterol	\rightarrow	Cholesterol
Na	→	Na
к	→	ĸ
Drug	\rightarrow	Drug
Na_to_K	\rightarrow	Na_to_K
View current fields O View OK Cancel	w unused field	settings

Figure 93. Edition du noeud Filtrer

Dans l'onglet Filtrer, cliquez sur les flèches situées en regard de *Na* et *K*. Des X rouges apparaissent sur les flèches pour indiquer que les champs sont désormais filtrés.

Reliez ensuite un noeud type connecté au noeud Filtrer. Le noeud type vous permet d'indiquer les types de champ utilisés, ainsi que la façon dont ils seront utilisés pour prédire les résultats.

Dans l'onglet Types, attribuez le rôle *Médicament* au champ **Cible** ; vous indiquez ainsi que le champ *Médicament* est celui sur lequel porte l'analyse. Laissez les autres champs paramétrés sur le rôle **Entrée** de sorte qu'ils soient utilisés comme prédicteurs.

Type	Annotations	Y		(
4- 000	🕶 🚺 🕨 Read Va	lues Clear '	Values	Clear All Value	s
Field -	Measurement	Values	Missing	Check	Role
📿 Age	🖉 Continuous	[15,74]		None	🔪 Input
A Sex	🖁 Flag	M/F		None	🔪 Input
A BP	Som Nominal	HIGH,LO		None	🔪 Input
A Cholesterol	🖁 Flag	NORMAL/		None	🔪 Input
A Drug	💑 Nominal	drugA,dru		None	🔪 Input 🚿
🛞 Na_to_K	🔗 Continuous	[6.268724		None	🔪 Input
					O Target
					Both
					None
View current	fields 🔘 View unus	ed field settings	;		Partition
					Solit
OK Cance					and spin
Cance					• • • • requenc
					Record ID

Figure 94. Edition du noeud type

Pour évaluer le modèle, placez un noeud C5.0 dans l'espace de travail et reliez-le à la fin du flux, comme l'indique l'illustration. Puis cliquez sur le bouton vert de la barre d'outil **Exécuter** pour exécuter le flux.



Figure 95. Ajout d'un noeud C5.0

Navigation dans le modèle



Figure 96. Navigation dans le modèle

Lorsque le noeud C5.0 est exécuté, le nugget de modèle est ajouté au flux et également à la palette Modèles en haut à droite de la fenêtre. Pour parcourir le modèle, cliquez avec le bouton droit de la souris sur une des icônes, puis sélectionnez **Modifier** ou **Parcourir** dans le menu contextuel.

Le navigateur de règles affiche l'ensemble des règles générées par le noeud C5.0 sous forme d'arbre décision. A l'ouverture du navigateur, l'arbre est réduit. Pour le développer et afficher tous les niveaux, cliquez sur le bouton **Tout**.

Model	Viewer	Summary	Settings	Annotation	s
b	14	1 2 3	B All	Ø•	18 (i)
∎- Na Na	a_to_K <= a_to_K >	: 14.64 14.64		<⇒ drug	Ŷ

Figure 97. Navigateur de règles

Ainsi, vous pouvez visualiser les éléments manquants. Pour les personnes ayant un rapport *Na*-sur-*K* inférieur à 14,64 et une pression artérielle élevée, l'âge détermine le choix du médicament. Pour les personnes présentant une faible pression artérielle, le taux de cholestérol semble être le prédicteur optimal.

Model	Viewer	Summary	Settings	Annotatio	ns
b	14	1 2 3		∳ •	1
P. Na	a_to_K <= ∙ BP = HI	: 14.64 GH			
	- Age	<= 50		🖙 dru	ıgA
	- Age	> 50		🖙 dru	βB
	BP = LC	WV			
	- Cho	lesterol = l	NORMAL		⊏> drugX
	- Cho	lesterol = I	HIGH		⇒ drugC
	BP = NO	ORMAL		🖙 dr	ugX
- Na	a_to_K ≻	14.64		🖙 dru	gY

Figure 98. Navigateur de règles développé au maximum

Vous pouvez consulter ce même arbre dans un format de graphique plus élaboré. Pour ce faire, cliquez sur l'onglet **Visualiseur**. Vous pouvez voir plus facilement le nombre d'observations contenues dans chaque catégorie de pression artérielle, ainsi que le pourcentage d'observations.



Figure 99. Arbre de décisions en format graphique

Utilisation d'un noeud Analyse

Vous pouvez évaluer l'exactitude du modèle à l'aide d'un noeud Analyse. Reliez un noeud Analyse (de la palette du noeud Sortie) au nugget de modèle, ouvrez le noeud Analyse et cliquez sur **Exécuter**.



Figure 100. Ajout d'un noeud Analyse

Le résultat du noeud Analyse indique que, avec ce jeu de données artificielles, le modèle a réalisé une prévision correcte du choix de médicament pour chaque enregistrement du jeu de données. Avec un jeu de données réel, il est peu probable que vous soyez confronté à une exactitude de 100 %. Vous pouvez néanmoins utiliser le noeud Analyse pour déterminer si le modèle a une exactitude acceptable pour votre application.

🔍 Analysi	s of [Drug]			
🐞 <u>F</u> ile	è Edit 🛛 🙀		14	
Analysis	Annotations	salet SV.		
8 Collaps	e All 🤤 E	xpand Al		
Results	for output field	Drug		
Con 🖻	paring \$C-Drug	with Dru	lg	
	Correct	200	100%	
	Wrong	0	0%	
8	Total	200		
				OK
				UK

Figure 101. Sortie du noeud Analyse

Chapitre 9. Filtrage des prédicteurs (sélection de fonction)

Le noeud Sélection de fonction vous permet d'identifier les champs les plus importants pour la prévision de certains résultats. A partir de centaines, voire de milliers de prédicteurs, le noeud Sélection de fonction filtre, classe et sélectionne celles qui peuvent être les plus importantes. En fin de compte, vous pouvez obtenir un modèle plus rapide et plus efficace (qui utilise moins de prédicteurs, s'exécute plus rapidement et est plus compréhensible).

Les données utilisées dans cet exemple représentent l'entrepôt de données d'un opérateur de téléphonie fictif et contiennent des informations concernant les réponses données par 5 000 clients de l'opérateur à une promotion spéciale. Ces données incluent un grand nombre de champs comprenant l'âge, la profession et les revenus des clients, ainsi que les statistiques d'utilisation de leur téléphone. Trois champs "cible" indiquent si le client a répondu à chacune des trois offres. L'opérateur souhaite utiliser ces données pour connaître les clients les plus susceptibles de répondre à des offres similaires à l'avenir.

Cet exemple utilise le flux nommé *featureselection.str*, qui fait référence au fichier de données nommé *customer_dbase.sav*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *featureselection.str* se trouve dans le répertoire des *flux*.

Cet exemple n'emploie comme cible que l'une des offres. Il utilise le noeud de création d'arbre CHAID pour développer un modèle visant à décrire les clients les plus susceptibles de répondre à la promotion. Il met en opposition deux approches :

- Sans la sélection de fonction. Tous les champs prédicteurs du jeu de données sont employés comme entrées pour l'arbre CHAID.
- Avec la sélection de fonction. Le noeud Sélection de fonction est utilisé pour sélectionner les 10 premiers prédicteurs. Ces prédicteurs sont ensuite employés comme entrées pour l'arbre CHAID.

Si nous comparons les deux modèles d'arbre obtenus, nous constatons que la sélection de fonction produit des résultats efficaces.

Création du flux



With All Fields

Figure 102. Exemple de flux Sélection de fonction

 Placez un noeud source Statistics sur un canevas de flux vide. Faites pointer ce noeud vers le fichier de données exemple *customer_dbase.sav*, disponible dans le répertoire *Demos* de votre installation IBM SPSS Modeler. (Vous pouvez également ouvrir le fichier de flux exemple *featureselection.str*, dans le répertoire des *flux*.) 2. Ajoutez un noeud type. Dans l'onglet Types, défilez vers le bas et modifiez le rôle du champ *response_01* en *Cible*. Modifiez le rôle en *Aucun* pour les autres champs de réponse (*response_02*) et *response_03*) ainsi que l'ID client (*custid* en haut de la liste. Laissez le rôle défini sur *Entrée* pour tous les autres champs, et cliquez sur le bouton **Lire les valeurs**, puis cliquez sur **OK**.

	view				0.1
ypes Format	Annotations				
- 00 0	🕨 🔰 🕨 Read Val	ues Clea	r Values	Clear All Va	alues
Field -	Measurement	Values	Missing	Check	Role
y ovvripua	Nominai	0,1		NONE	a input
ownpc 🧉	nominal	0,1		None	🔪 Input
🕻 ownipod	📩 Nominal	0,1		None	🔪 Input
🕻 owngame 🧯	📩 Nominal	0,1		None	🔪 Input
ownfax 🤞	📩 Nominal	0,1		None	🔪 Input
news	Nominal	0,1		None	tuqni 🖌
response 01	Nominal	0.1		None	O Target
response 02	Nominal	0.1		None	O None
response_03	Nominal	0,1		None	○ None
View current fi	elds 🔘 View unus	ed field settin	as		
s non carron n		ou noid couir	.90		

Figure 103. Ajout d'un noeud type

- **3**. Ajoutez au flux un noeud de modélisation Sélection de fonction. Dans ce noeud, vous pouvez définir les règles et les critères de filtrage ou de désactivation des champs.
- 4. Exécutez le flux afin de créer le nugget de modèle Sélection de fonction.
- 5. Faites un clic droit sur le nugget de modèle dans le flux ou dans la palette Modèles et choisissez **Modifier** ou **Parcourir** pour afficher les résultats.

😡 resp	onse_01							X
-	📦 File	🏷 Generate	Pr	eview			0 -	
						-		
Model	Summary	Annotations						
	- (e)							
		Rank	- 4	<u>1</u>				
	Repk	Field		Maacuramant		Importance	Valua	
	TXGITIX 4	A ed		optipuous	-	Importance	1.0	
	2		A N	ominal	÷	Important	1.0	
	4	ownpc dedcat		rdinal	÷	Important	1.0	-
	4		N.	ominal		Important	1.0	
	4		A N	ominal		Important	1.0	
	6			ominal		Important	1.0	
	7	equinmon	AC	ontinuous		Important	1.0	
	. 8		AN	ominal	-	Important	1.0	
	q	Č ehill	A N	ominal	-	Important	1.0	
	10		A N	ominal	-	Important	1.0	
	11	of forward	A N	ominal		Important	1.0	
	12	tollmon	1 C	ontinuous		Important	1.0	
	13		A N	ominal	-	Important	1.0	
	14		A N	ominal	-	Important	1.0	
	15		A N	ominal		Important	1.0	
	16	equipten	AC	ontinuous	*	Important	1.0	
	17	otollfree	A N	ominal	*	Important	1.0	
	18	toliten	AC	ontinuous	*	Important	1.0	
	19	Č churn	A N	ominal	*	Important	1.0	
	20	🛆 snousedcat	10	rdinal	+	Important	1.0	-
Selecte	ed fields: 34	Total fields availa	able: 12 95 🕂	8 <= 0.95 💽 < 0.	9			
-			~ ~					
			9 Scre	ened Fields				
	Field 🐬	Measureme	nt		R	eason		
	父 ownyci	r 💑 Nominal		Single category to	o la	rge		-
	父 owntv	nominal ស		Single category to	o la	rge		
	📿 owndv	d 🛛 💑 Nominal		Single category to	o la	rge		
	📿 owned	nominal 🔬		Single category to	o la	rge		
	🛞 Inwirete	en 🔗 Continuous		Too many missing	val	ues		
	🛞 Inwiren	n 🔗 Continuous		Too many missing	val	ues		
	🛞 Inequip	🖉 Continuous		Coefficient of vari	atio	n below thre	shold	

Figure 104. Onglet Modèle dans le nugget de modèle Sélection de fonction

Le panneau supérieur contient les champs considérés comme utiles pour la prévision. Ils sont classés en fonction de leur importance. Le panneau inférieur indique les champs filtrés, et la raison du filtrage. En examinant les champs du panneau supérieur, vous pouvez décider de ceux à utiliser lors des sessions de modélisation suivantes.

- 6. Nous pouvons désormais sélectionnez les champs à employer en aval. Bien que 34 champs aient été identifiés à l'origine comme étant importants, nous souhaitons quand même réduire davantage l'ensemble de prédicteurs.
- 7. Sélectionnez uniquement les 10 premières prédicteurs en décochant les cases de la première colonne pour désélectionner les prédicteurs superflus. (Cliquez sur la coche de la ligne 11, maintenez la touche Maj. appuyée et cliquez sur la coche de la ligne 34). Fermez le nugget de modèle.
- 8. Pour comparer les résultats sans sélection de fonction, ajoutez deux noeuds de modélisation CHAID au flux : un noeud utilisant la sélection de fonction et un noeud ne s'en servant pas.
- 9. Connectez un noeud CHAID au noeud type et l'autre au nugget de modèle Sélection de fonction.
- **10.** Ouvrez chacun des noeuds CHAID, sélectionnez l'onglet Options de création et vérifiez que les options **Créer un nouveau modèle**, **Créer un seul arbre** et **Lancer une session interactive** sont sélectionnées dans le panneau Objectifs.

Dans le panneau Options de base, vérifiez que la Profondeur maximale de l'arbre est définie sur 5.



Figure 105. Paramètres des objectifs du noeud de modélisation CHAID pour tous les champs prédicteurs

Création des modèles

- 1. Exécutez le noeud CHAID qui emploie tous les prédicteurs du jeu de données (celui connecté au noeud type). Notez la durée du traitement. La fenêtre de résultats affiche un tableau.
- 2. Dans les menus, choisissez Arbre > Développer l'arbre pour développer et afficher l'arbre.

The active Tree of CHAID #6	_ 🗆 🛛
🙀 File 📄 Edit 🖋 View Iree 🖏 Generate 🕘 🚮	@ ×
Viewer Gains Risks Annotations	
🕒 ြ 🎧 ि 🕘 🔳 🖷 🙊 🙊 🕮 🕮 🏦 🤹 🎥 🛞	
regnance 01	
Category % n	
■ 0.000 91.640 4582 ■ 1.000 8.360 418	
Total 100.000 5000	
ownpc	
Au). P-value=0.000, Chi-square=57.452, u=1	
	+
🛄 😨 🐮 💽 😴 🧐 💸	Data
	ОК

Figure 106. Développement de l'arbre dans le Générateur d'arbres

3. A présent, procédez de même pour l'autre noeud CHAID, qui n'utilise que 10 prédicteurs. Là encore, développez l'arbre lorsque le Générateur d'arbres s'ouvre.

Le second modèle s'exécute normalement plus vite que le premier. Le jeu de données considéré étant relativement petit, la différence de temps d'exécution est peut-être de quelques secondes seulement, mais, pour des jeux de données réels, plus volumineux, cette différence peut s'avérer très importante (plusieurs minutes, voire plusieurs heures). Utiliser la sélection de fonction peut accélérer considérablement vos temps de traitement.

Le second arbre contient également moins de noeuds d'arbre que le premier. Il est plus simple à comprendre. Toutefois, avant de décider de vous en servir, vous devez vérifier qu'il est efficace et le comparer au modèle utilisant tous les prédicteurs.

Comparaison des résultats

Pour comparer les deux résultats, nous devons utiliser une mesure d'efficacité. Pour ce faire, utilisons l'onglet Gains du Générateur d'arbres. Portons notre attention sur le **lift**, qui mesure le degré de probabilité selon lequel les enregistrements d'un noeud peuvent faire partie de la catégorie cible, comparés à tous les enregistrements du jeu de données. Par exemple, une valeur de lift (augmentation) de 148 % signifie que les enregistrements du noeud ont 1,48 fois plus de chances d'appartenir à la catégorie cible que tous les enregistrements du jeu de données. Le lift est spécifié dans la colonne *Index* de l'onglet Gains.

- Dans le Générateur d'arbres de l'ensemble complet des prédicteurs, cliquez sur l'onglet Gains. Définissez la catégorie cible sur 1,0. Passez à un affichage en quartiles. Pour ce faire, cliquez d'abord sur le bouton de la barre d'outils Quantiles. Puis sélectionnez Quartile dans la liste déroulante à droite de ce bouton.
- 2. Répétez cette procédure dans le Générateur d'arbres pour l'ensemble des 10 prédicteurs, de sorte à avoir deux tableaux de gains similaires à comparer, comme l'illustrent les figures suivantes.

🔰 <u>F</u> ile 🛛 📄 Ed	it 🥑 View	v <u>T</u> ree	🕙 Generate 🛛 🚺					0
/iewer Gains	Risks Anr	notations						
2 1 1	Quartile		▼ 🖽 🖉 Ga	ins	👻 🍾 Targ	get category 1.0	• I	7
			Target va	riable: response	01 Target catego	prv: 1.0		
aining Sample								
odes		Percentile	Percentile: n	Gain: n	Gain (%)	Response (%)	Index (%)	
29,43,8,42,38,	53,45,49,33	25.00	1250.00	231.00	55.29	18.49	221.17	
56,21,22,62,59	41,40,51,	50.00	2500.00	358.00	85.54	14.30	171.09	
47,32,55,58,19	,46	75.00	3750.00	407.00	97.45	10.86	129.94	
23,52,60,37,50	39,35,57,	100.00	5000.00	418.00	100.00	8.36	100.00	
1				Anadoseden				
Interactive	Tree of CH	AID #1					_	
Eila 🕒 Ed	it 🥑 View	v <u>T</u> ree	🏷 <u>G</u> enerate 🛛 🚺					
/iewer Gains	Risks Anr	notations						
/iewer Gains	Risks Ann	notations	▼ Ⅲ ∠ Ga	ins	Taro	aet category 1.0		7
/iewer Gains	Risks Anr Quartile	notations	▼ Ⅲ ∠ Ga	ins	👻 🎉 Targ	get category 1.0	•	ý
/iewer Gains	Risks Anr	notations	▼ Ⅲ ∠ Ga Taroet va	iins riable: response	Target catego	get category 1.0	•	
/iewer Gains	Risks Ann	notations	Target va	iins riable: response_	01 Target catego	yet category 1.0 pry: 1.0	•	,
/iewer Gains	Risks Anr	notations	Target va	iins riable: response_	▼ Starg 01 Target catego	yet category <mark>1.0</mark> vry: 1.0	•	7
Viewer Gains	Risks Anr Quartile	notations	Target va	ins riable: response_ Gain: n	▼ Starg 01 Target catego Cain (%)	get category 1.0 pry: 1.0) Response (%)	• Index (%)	
Aning Sample	Risks Ann Quartile Percenti 25.00	Inctations	Target va	ins riable: response_ Gain: n 203.00	01 Target catego Gain (%) 48.45	get category 1.0 yry: 1.0) Response (%) 16.20	Index (%)	
A construction of the second s	Risks Ann Quartile Percenti 25.00 50.00	Ile		iins riable: response_ Gain: n 203.00 308.00	01 Target catego Gain (%) 48.45 73.57	yet category 1.0 pry: 1.0) Response (%) 16.20 12.30	Index (%) 193.81 147.14	
A construction of the cons	Risks Ann Quartile Percenti 25.00 50.00 75.00	ile	Target va Target va Percentile: n 1250.00 2500.00 3750.00	ins riable: response_ Gain: n 203.00 308.00 385.00	C1 Target catego Gain (%) 48.45 73.57 92.14	yet category 1.0 rry: 1.0 Response (%) 16.20 12.30 10.27	Index (%) 193.81 147.14 122.86	
aining Sample	Risks Ann Quartile Percentil 25.00 50.00 75.00 100.00	ile	E C Ga Target va Percentile: n 1250.00 2500.00 3750.00 5000.00	riable: response_ Gain: n 203.00 308.00 385.00 418.00	Cain (%) Gain (%) 48.45 73.57 92.14 100.00	yet category 1.0 pry: 1.0 () Response (%) 16.20 12.30 10.27 8.36	Index (%) 193.81 147.14 122.86 100.00	
Aining Sample odes 323,15,12 1,26,10,7 17,11,20 1,24,16,19,25	Risks Ann Quartile Percentil 25.00 50.00 75.00 100.00	ile	Target va	ins riable: response_ 203.00 308.00 385.00 418.00	Cain (%) Gain (%) 48.45 73.57 92.14 100.00	yet category 1.0 pry: 1.0 Response (%) 16.20 12.30 10.27 8.36	Index (%) 193.81 147.14 122.86 100.00	

Figure 107. Graphiques de gain des deux modèles CHAID

Chaque tableau de gains regroupe les noeuds terminaux de son arbre en quartiles. Pour comparer l'efficacité des deux modèles, examinez le lift (valeur *Index*) du quartile supérieur dans chaque tableau.

Si tous les prédicteurs sont inclus, le modèle affiche un lift (augmentation) de 221 %. Plus précisément, les observations présentant les caractéristiques de ces noeuds ont 2,2 fois plus de chances de répondre à la promotion cible. Pour connaître ces caractéristiques, cliquez sur la ligne supérieure afin de la sélectionner. Passez ensuite à l'onglet Visualiseur, où les noeuds correspondants sont désormais mis en évidence en noir. Parcourez l'arbre de haut en bas, jusqu'à chaque noeud terminal mis en évidence, afin de voir comment les prédicteurs ont été divisés. A lui seul, le quartile supérieur comprend 10 noeuds. Convertis en modèles de scoring réels, 10 profils client différents peuvent être difficiles à gérer.

Avec l'inclusion des 10 premiers prédicteurs (identifiés par la sélection de fonction) seulement, le lift (augmentation) est de presque 194 %. Bien que ce modèle ne soit pas aussi performant que celui employant tous les prédicteurs, il est indéniablement utile. Dans ce cas, le quartile supérieur n'inclut que quatre noeuds et est donc plus simple. Nous arrivons par conséquent à la conclusion qu'il est préférable d'utiliser le modèle Sélection de fonction au lieu de celui employant tous les prédicteurs.

Récapitulatif

Passons à présent en revue les avantages de la sélection de fonction. Utiliser moins de prédicteurs est plus économique. En effet, vous avez moins de données à collecter, à traiter et à intégrer dans vos modèles. Le temps de calcul s'en trouve amélioré. Dans cet exemple, même avec l'étape supplémentaire de la sélection de fonction, la création du modèle a été nettement plus rapide avec l'ensemble réduit de prédicteurs. Avec un jeu de données réel plus volumineux, les gains de temps seraient considérables.

Utiliser moins de prédicteurs simplifie le scoring. Comme le montre cet exemple, vous ne pouvez identifier que quatre profils de clients susceptibles de répondre à la promotion. Veuillez noter qu'avec des quantités plus importantes de prédicteurs, vous risqueriez de surajuster votre modèle. Il est possible que le modèle le plus simple se généralise mieux aux autres jeux de données (mieux vaut néanmoins effectuer un test à titre de vérification).

Pour la sélection de fonction, vous auriez pu utiliser un algorithme de création d'arbre. L'arbre identifie ainsi automatiquement les prédicteurs les plus importants. En fait, l'algorithme CHAID est souvent utilisé à cet effet et il est même possible de développer l'arbre niveau par niveau pour en contrôler la profondeur et la complexité. Toutefois, le noeud Sélection de fonction est plus rapide et plus facile à utiliser. Il classe tous les prédicteurs en une seule fois, ce qui vous permet d'identifier rapidement les champs les plus importants. En outre, il vous offre la possibilité de changer le nombre de prédicteurs à inclure. Vous pouvez facilement réappliquer cet exemple, en utilisant cette fois les 15 ou 20 premiers prédicteurs au lieu des 10 premiers, afin de comparer les résultats et de déterminer le modèle optimal.
Chapitre 10. Réduction de la longueur des chaînes de données d'entrée (Noeud Recoder)

Réduction de la longueur des chaînes de données d'entrée (Reclassifier)

Pour les modèles de régression logistique et de Discriminant automatique qui incluent un modèle de régression logistique binomiale, les champs de type chaîne sont limités à 8 caractères maximum. Lorsque les chaînes contiennent plus de 8 caractères, elles peuvent être recodées à l'aide du noeud Recoder.

Cet exemple utilise le flux intitulé *reclassify_strings.str*, qui référence le fichier de données *drug_long_name*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *reclassify_strings.str* se trouve dans le répertoire des *flux*.

Cet exemple se concentre sur une petite partie d'un flux et présente le type d'erreurs pouvant être générées avec des chaînes trop longues et explique comment utiliser le noeud Recoder pour modifier les détails des chaînes et leur donner une longueur acceptable. Bien que cet exemple utilise un noeud de régression logistique binomiale, il convient également lors de l'utilisation du noeud Discriminant automatique pour générer un modèle de régression logistique binomiale.

Reclassification des données

1. A l'aide d'un noeud source Délimité, connectez-vous au jeu de données *drug_long_name* dans le dossier *Demos*.



Figure 108. Flux d'échantillons présentant une reclassification de chaînes pour une régression logistique binomiale

- 2. Ajoutez un noeud type au noeud source et sélectionnez Cholesterol_long comme cible.
- 3. Ajoutez un noeud Régression logistique au noeud type.
- 4. Dans le noeud Régression logistique, cliquez sur l'onglet Modèle et sélectionnez la procédure **Binomial**.

	A				
	Read Value	ues Clear \	/alues	Clear All Valu	es
Field -	Measurement	Values	Missing	Check	Role
Age	🔗 Continuous	[15,74]	_	None	🔪 Input
Sex	🎖 Flag	M/F		None	🔪 Input
BP	💑 Nominal	HIGH,LO		None	🔪 Input
Na	🔗 Continuous	[0.500517		None	🔪 Input
×к	🔗 Continuous	[0.020152		None	🔪 Input
Drug	💑 Nominal	drugA,dru		None	🔪 Input
Cholesterol	🎖 Flag	"Normal le		None	🔘 Target

Figure 109. Détails de chaînes de grande longueur dans le champ "Cholesterol_long"

5. Lorsque vous exécutez le noeud Régression logistique dans *reclassify_strings.str*, un message d'erreur apparaît pour vous prévenir que les valeurs de chaîne **Cholesterol_long** sont trop longues.

Si vous rencontrez ce genre de messages d'erreur, suivez la procédure expliquée dans le reste de cet exemple pour modifier vos données.

Message	
(i) Stream execution started	
😵 Field 'Cholesterol_long' has value 'High level of cholesterol' that is too long.	
😵 Field 'Cholesterol_long' has value 'Normal level of cholesterol' that is too long.	
 Stream execution complete, Elapsed=0.39 sec, CPU=0.02 sec 	
A Execution was interrupted	

Figure 110. Message d'erreur affiché lors de l'exécution du noeud de régression logistique binomiale

- 6. Ajoutez un noeud Recoder au noeud type.
- 7. Dans le champ Recoder, sélectionnez Cholesterol_long.
- 8. Saisissez Cholesterol comme nouveau nom de champ.
- 9. Cliquez sur le bouton **Obtenir** pour ajouter les valeurs**Cholesterol_long** à la colonne de valeurs d'origine.
- 10. Dans la nouvelle colonne de valeurs, saisissez **Elevé** à côté de la valeur d'origine du **Niveau de cholestérol élevé** et **Normal** à côté de la valeur d'origine de **Niveau de cholestérol normal**.

Cholesterol		×
Preview		0
Settings Annotations		
Mode:	💿 Single 🔘 Multiple	
Reclassify into:	New field C Existing field	1
Reclassify field:		
Cholesterol_long		-
New field name:		
Cholesterol		
Reclassify values:	Clear new	4 Auto
Original value	New value	
High level of cholesterol	High	
Normal level of cholesterol	Normal	
		÷
For unspecified values use: O	riginal value 🔘 Default value	undef
OK Cancel		Apply Reset

Figure 111. Reclassification des chaînes longues

- 11. Ajoutez un noeud Filtrer au noeud Recoder.
- 12. Dans la colonne Filtrer, cliquez pour supprimer Cholesterol_long.

Filter		
Filter Annotations		
7.	Fields	: 8 in, 1 filtered, 0 renamed, 7 out
Field -	Filter	Field
Age	\rightarrow	Age
Sex	\rightarrow	Sex
BP	\rightarrow	BP
Na	\rightarrow	Na
к	\rightarrow	к
Drug	\rightarrow	Drug
Cholesterol_long	×>	Cholesterol_long
Cholesterol	\rightarrow	Cholesterol
View current fields View	v unused field s	ettings

Figure 112. Filtrage du champ "Cholesterol_long" à partir des données

13. Ajoutez un noeud type au noeud Filtrer et sélectionnez Cholesterol comme cible.

Type	wiew			0	
	Read Valu	ies Clear V	alues	Clear All Values	
Field -	Measurement	Values	Missing	Check	Role
🔆 Age	🖉 Continuous	[15,74]		None	🔪 Input
A Sex	🖁 Flag	MÆ		None	🔪 Input
A BP	Nominal	HIGH,LO		None	🔪 Input
🚯 Na 🛛 🐰	🖉 Continuous	[0.500517		None	🔪 Input
ŷк .	Continuous	[0.020152		None	🔪 Input
A Drug	🇞 Nominal	drugA,dru		None	🔪 Input
A Cholesterol	🎖 Flag	Normal/High		None	🔘 Target
View current fi Cancel	ields 🔘 View unus	ed field settings		A	pply Res

Figure 113. Détails de chaînes courtes dans le champ "Cholesterol"

- 14. Ajoutez un noeud Logistique au noeud type.
- 15. Dans le noeud Logistique, cliquez sur l'onglet Modèle et sélectionnez la procédure Binomial.
- **16**. Vous pouvez maintenant exécuter le noeud Logistique binomiale et générer un modèle sans qu'un message d'erreur ne s'affiche.

😡 Cholesterol			×
Fields Model Expert	Analyze Annotations		
Model name: 💿 Auto 🔘	Custom		
👿 Use partitioned data			
👿 Build model for each sp	lit		
Procedure: O Multinom	ial	Binomial	
Binomial Procedure			
Method: Enter			
Categorical Inputs:			
Field Name	Contrast	Base Category	
		×	
🗹 Include constant in equa	ation		
OK 🕨 Run Cance	el	Apply	≡t

Figure 114. Choix de la procédure Binomial

Cet exemple ne présente qu'une partie d'un flux. Si vous avez besoin d'informations supplémentaires sur les types de flux dans lesquels vous pouvez avoir besoin de reclassifier de longues chaînes, les exemples suivants sont disponibles :

- Noeud Discriminant automatique. Pour plus d'informations, voir la rubrique «Modélisation de la réponse client (Discriminant automatique)», à la page 39.
- Noeud Régression logistique binomiale. Pour plus d'informations, voir la rubrique Chapitre 13, «Attrition dans le domaine des télécommunications (régression logistique binomiale)», à la page 141.

Des informations supplémentaires sur l'utilisation d'IBM SPSS Modeler, telles que le guide de l'utilisateur, le guide de référence des noeuds et le guide des algorithmes, sont disponibles dans le répertoire *Documentation* du disque d'installation.

Chapitre 11. Modélisation de la réponse client (Liste de décision)

L'algorithme Liste de décision génère des règles qui indiquent une probabilité plus ou moins élevée d'obtenir un résultat binaire (oui ou non) donné. Les modèles Liste de décision sont largement utilisés dans la gestion de la relation client, par exemple dans les centres d'appel ou les applications marketing.

Cet exemple repose sur une société fictive qui souhaite obtenir des résultats plus rentables au cours des prochaines campagnes de marketing en présentant à chaque client une offre adaptée. En particulier, l'exemple utilise un modèle Liste de décision pour identifier les caractéristiques des clients les plus à même de répondre favorablement, sur la base des promotions précédentes, et de générer un fichier d'adresses en fonction des résultats.

Les modèles Liste de décisions sont particulièrement adaptés à la modélisation interactive et vous permettent de régler les paramètres du modèle et d'obtenir des résultats immédiats. Si vous souhaitez utiliser une autre approche qui vous permet de créer automatiquement plusieurs modèles différents et de classer les résultats obtenus, utilisez le noeud Discriminant automatique.



Figure 115. Flux d'échantillons Liste de décision

Cet exemple utilise le flux *pm_decisionlist.str*, qui fait référence au fichier de données *pm_customer_train1.sav*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *pm_decisionlist.str* se trouve dans le répertoire des *flux*.

Données d'historique

Le fichier *pm_customer_train1.sav* comporte des données d'historique suivant les offres faites à des clients spécifiques au cours de campagnes passées, comme l'indique la valeur du champ *campaign*. Le plus grand nombre d'enregistrements se trouve dans la campagne *Premium account*.

🝃 <u>F</u> ile	📄 <u>E</u> dit (🖞 <u>G</u> enerate 🛛 🚺					0
Table	Annotations				_	_	
	customer_id	campaign	response	response_date	purchase	purchase_date	product_id
1	7	Premium account	0	\$null\$	0	\$null\$	\$null\$
2	13	Premium account	0	\$null\$	0	\$null\$	\$null\$
3	15	Premium account	0	\$null\$	0	\$null\$	\$null\$
4	16	Premium account	1	2006-07-05 00:00:00	0	\$null\$	183
5	23	Premium account	0	\$null\$	0	\$null\$	\$null\$
6	24	Premium account	0	\$null\$	0	\$null\$	\$null\$
7	30	Premium account	0	\$null\$	0	\$nuli\$	\$null\$
8	30	Gold card	0	\$null\$	0	\$null\$	\$null\$
9	33	Premium account	0	\$null\$	0	\$null\$	\$null\$
10	42	Gold card	0	\$null\$	0	\$null\$	\$null\$
11	42	Premium account	0	\$null\$	0	\$null\$	\$null\$
12	52	Premium account	0	\$null\$	0	\$null\$	\$null\$
13	57	Premium account	0	\$null\$	0	\$null\$	\$null\$
14	63	Premium account	1	2006-07-14 00:00:00	0	\$null\$	183
15	74	Premium account	0	\$null\$	0	\$null\$	\$null\$
16	74	Gold card	0	\$null\$	0	\$null\$	\$null\$
17	75	Premium account	0	\$null\$	0	\$null\$	\$null\$
18	82	Premium account	0	\$null\$	0	\$null\$	\$null\$
19	89	Gold card	0	\$null\$	0	\$null\$	\$null\$
20	89	Premium account	0	\$null\$	0	\$null\$	\$null\$
	4			And and the second s		dented and the second	1

Figure 116. Données sur les anciennes promotions

Les valeurs du champ *campaign* sont en fait codées comme des entiers dans les données, avec des libellés définies dans le noeud type (par exemple 2 = *Premium account*). Vous pouvez masquer ou afficher les libellés de valeur dans le tableau à l'aide de la barre d'outils.

Le fichier inclut aussi un certain nombre de champs contenant des informations démographiques et financières sur chaque client qui peut servir à créer ou à "former" un modèle qui prévoit les taux de réponse pour différents groupes en fonction de caractéristiques spécifiques.

Création du flux

1. Ajoutez un noeud source Statistics qui pointe sur *pm_customer_train1.sav*, dans le dossier *Demos* du répertoire d'installation d'IBM SPSS Modeler. (Vous pouvez spécifier \$CLE0_DEMOS/ dans le chemin du fichier comme raccourci de référence de ce dossier)



Figure 117. Lecture de données

2. Ajoutez un noeud type, puis sélectionnez *Réponse* en tant que champ cible (Rôle = **Cible**). Paramétrez le niveau de mesure de ce champ sur **Indicateur**.

Type Pre	Annotations				× • • • • •
√ - ∞ ∞	 Read Val 	lues Clear	Values	Clear All Va	alues
Field -	Measurement	Values	Missing	Check	Role
🚫 customer id 🌡	🔗 Continuous	[7,116993]		None	🛇 None 🛛 🖊
🚫 campaign 🛛	Nominal	1,2,3,4		None	O None
🚫 response	🎖 Flag	1/0		None	🔘 Target
response 🖌	Continuous	[2006-04		None	O None
🚫 purchase 🛛 🤞	🖉 Continuous	[0,1]		None	
purchase 🖌	🖉 Continuous	[2006-04		None	
🚫 product_id 💡	🖉 Continuous	[183,421]		None	None None
🚫 Rowid 🛛 🤞	🔗 Continuous	[1,19599]		None	🛇 None 📃
ana 🛆	🖉 Continuous	M0.061		None	🔪 Innut 💽
View current f OK Cancel	ields 🔘 View unu:	sed field setting	3 8		Apply Reset

Figure 118. Configuration du niveau de mesure et du rôle

- **3**. Définissez l'option role sur **None** (aucun) pour les champs suivants : *customer_id, campaign, response_date, purchase, purchase_date, product_id, Rowid* et X_*random*. Ces champs ont tous des utilisations dans les données mais ne seront pas utilisés pour la création du modèle réel.
- 4. Cliquez sur le bouton Lire les valeurs dans le noeud type pour vérifier que les valeurs sont instanciées.

Les données incluent des informations sur quatre campagnes différentes, mais vous vous concentrerez sur l'analyse d'une seule campagne à la fois. Comme le plus grand nombre d'enregistrements se trouve dans

la campagne Premium (codée *campaign=2* dans les données), vous pouvez utiliser un noeud Sélectionner pour n'inclure que ces enregistrements dans le flux.

Select		×
-?>	Preview	
Settings	Annotations	
Mode:	💿 Include 🔘 Discard	
Condition:	campaign = 2	
ОК Саг	cel	Apply

Figure 119. Sélection d'enregistrements pour une seule campagne

Création du modèle

1. Reliez un noeud Liste de décision au flux. Dans l'onglet Modèle, définissez la **Valeur cible** sur 1 pour indiquer le résultat que vous souhaitez rechercher. Dans notre exemple, vous recherchez des clients qui ont répondu *Oui* à une offre précédente.

🚱 response[1]	
Fields Model Expert Analyze Annotations	3
Model name: 💿 Auto 🛇 Cust	om
👿 Use partitioned data	
Suild model for each split	
Mode: O Generate model O Launch interacti	ve session nation
Target value:	1
Find segments with: Maximum number of segments:	High probability 🐨
Minimum segment size As percentage of previous segment (%): As absolute value (N):	5.0 * 50 *
Segment rules Maximum number of attributes: Allow attribute re-use Confidence interval for new conditions (%):	5 - 85.0 -
OK Run Cancel	Apply Reset

Figure 120. Noeud Liste de décision, onglet Modèle

- 2. Sélectionnez Lancer une session interactive.
- **3**. Pour conserver la simplicité du modèle pour cet exemple, paramétrez le nombre maximum de segments sur 3.
- 4. Changez l'intervalle de confiance pour les nouvelles conditions à 85 %.
- 5. Dans l'onglet Expert, définissez le Mode sur Expert.

😡 response[1]		
Fields Model Expert Analyze	Annotations	0
Mode: 🔘 Simple 💿 Expert		
Binning method:	Equal Count 🔝	
Number of bins:	10 ≑	
Model search width:	5 ≑	
Rule search width:	5 ≑	
Bin merging factor:	2.0 ≑	
Allow missing values in conditi	ons	
V Discard intermediate results		
Interactive mode only		
Maximum number of alternatives:	3 🗧	
OK 🕨 Run Cancel		Apply Reset

Figure 121. Noeud Liste de décision, onglet Expert

- 6. Augmentez le **Nombre maximal d'alternatives** à 3. Cette option fonctionne en association avec le paramètre **Lancer une session interactive** que vous avez sélectionné dans l'onglet Modèle.
- 7. Cliquez sur Exécuter pour afficher le visualiseur Liste interactive.

🍮 Inte	ractive List: response[1] #1					
🐞 File	📄 Edit 💰 View Iools 🖔 Generate 🛛 🗐	<u>-</u>	8) 🕜 📼 🐔 😫	-	0	X
Viewer	Gains Annotations					
- •	ake Snapshot		Segment Finder	High Probability 💌	Settings	
Targe	t field: 🔎 response		Max. no. of new segment	s: 3 ⁴	Find Segments	
Targe	t value: 1					
id	Segment Rules	Score	Cover (n)	Frequency	Probability	*
	All segments including Remainder		13,504	1,952	2 14.45%	1
	Remainder		13,504	1,952	2 14.45%	
Model	Summary; Cover 0: Frequency 0: Probability 0%					
						ОК

Figure 122. visualiseur Liste interactive

Etant donné qu'aucun segment n'a encore été défini, tous les enregistrements sont inclus dans le reste. Sur les 13 504 enregistrements que compte l'échantillon, 1 952 ont dit *Oui*, soit un taux de correspondance global de 14,45%. Vous souhaitez améliorer ce taux en identifiant les segments de clients les plus (ou les moins) susceptibles de donner une réponse favorable.

8. Dans les menus du visualiseur Liste interactive, choisissez :

Outils > **Rechercher les segments**

Interactive List: response	se[1] #1					
🐞 File 🌛 Edit 💰 View	Tools 🖔 Generate 🛛 💼		🚺 🖏 🕅 🖬 🍫		0	×
Viewer Gains Annotations Viewer Gains Annotations Image: Target field: OP response Target value: 1	Find Segments Settings Organize Model Measures Organize Data Selections Change Target Value Take Snapshot		Segment Finder Find segments with: Max.no.of new segments	High Probability 🔻	Settings	
id Segment Rules		Score	Cover (n)	Frequency	Probability	
All segments including Re	emainder		13,504	1,952	14.45%	4
Model Summary; Cover 0: Frequ	iency 0: Probability 0%					ナ + + × に ・ ・ ・ ・ ・ ・ ・ ・ ・
					(ок

Figure 123. Visualiseur Liste interactive

Cette option exécute la tâche d'exploration par défaut sur la base des paramètres que vous avez définis dans le noeud Liste de décision. La tâche terminée renvoie trois modèles alternatifs, qui sont répertoriés dans l'onglet Alternatives de la boîte de dialogue Albums de modèles.

Name	T	arget	N	o. of Segme	nts	Cover		Freq.	Prob.
Alternative 1	1				3	-	2,375	5 1,348	56.76
Alternative 2	1				3		2,368	3 1,326	56.00
Alternative 3	1				3		2,380) 1,329	55.84
Alternative Pr	eview								
id	Segment	Rules			Score		Cover (n)	Frequency	Probability
	All segm	ents including	, Remainde	r	-		13,504	1,952	14.459
1	⊡ inco inc nu	me, numbe :ome > 55267 mber_produc	r_produc 7.000 and cts > 1.000	ts	1		912	795	87.179
2	⊡ rfm_ rfr nu	_score, num n_score > 12 mber_transa	nber_tran ∴333 and ctions > 2.	sactions	1		737	360	48.859
3	num nu	ber_transa mber_transa mber_transa :ome > 46072	ctions, in ctions > 0. ctions <= 1 2.000	come DOO and .000 and	1		731	174	23.809
	Remaind	er					11,124	623	5.60%
	Snapshots	3							

Figure 124. Modèles alternatifs disponibles

9. Sélectionnez la première alternative dans la liste ; ses détails sont affichés dans le panneau Aperçu de l'alternative.

Name		Target	No. of S	eqments	Cover		Freq.	rob.
Alternative 1		1.0		3		2,375	1,348	56.76%
Alternative 2		1.0		3		2,368	1,326	56.00%
Alternative 3		1.0		3		2,380	1,329	55.84%
Alternative I	Preview							
id	Segme	ent Rules		Scor	e	Cover (n)	Frequency	Probability
	All seg	ments including	Remainder			13,504	1,952	14.45%
1	□ income, number_products 1 income > 55267.000 and number_products > 1.000			1.0		912	2 795	87.17%
2	⊟ rfr	n_score, num rfm_score > 10. number_transad	ber_transacti 535 and :tions > 3.000	ons 1.0		725	; 357	49.24%
3	e av	erage#balance average#balanc average#balanc number_produc rfm_score > 9.2	e #feed#index, :e#feed#index > :e#feed#index < ts <= 2.000 and 39	numb∢ > 0.000 { <= 349.01.0		738) 196	26.56%
	Remair	nder				11,129	604	5.43%
▲ Load	Snapsh	ots						

Figure 125. Modèle alternatif sélectionné

Le panneau Aperçu de l'alternative vous permet de parcourir rapidement plusieurs alternatives sans changer le modèle de travail, ce qui facilite l'expérimentation de différentes approches.

Remarque : Pour mieux voir le modèle, vous pouvez agrandir le panneau Aperçu de l'alternative dans la boîte de dialogue comme l'indique l'illustration. Pour ce faire, faites glisser la bordure du panneau.

En utilisant des règles basées sur les prédicteurs comme le revenu, le nombre de transactions par mois et le score RFM, le modèle identifie des segments avec des taux de réponse qui sont plus élevés que ceux de l'ensemble de l'échantillon. Lorsque les segments sont combinés, ce modèle suggère que vous pouvez améliorer votre taux de correspondance jusqu'à 56,76%. Cependant, le modèle ne couvre qu'une petite portion de l'échantillon global et plus de 11 000 enregistrements (dont plusieurs centaines de correspondances) sont inclus dans le reste. Vous recherchez un modèle qui capture davantage de ces correspondances tout en excluant toujours les segments peu performants.

10. Pour essayer une autre approche de modélisation, sélectionnez les options suivantes dans les menus :

Outils > Paramètres

Create/Edit Mining Task: resp	onse[1]		X
Load Settings: response[1]	•	New 🗙	
🎯 Target Field: 💧	🔎 response	Target Value: 1	
Simple Settings			_
Find segments with:		High Probability 🤝	
Maximum number of new segments	:	3	
Minimum segment size			
As percentage of previous segn	nent (%):	5.0 ≑	
As absolute value (N):		50 🚔	
Maximum number of alternatives:		3	
Maximum attributes per segment:		5	
Allow attribute re-use within	segment		
Confidence interval for new condition	ons (%):	85.0	
Expert Settings			_
Binning method:	Equal Count	Number of bins:	10
Model search width:	5	Rule search width:	5
Bin merging factor:	2.00		
Allow missing values in conditions:	True	Discard intermediate results:	True
Data			
Build Selection: All Data	▼ 🖂		
Available fields: 🔘 All fields 🔘 C	Custom Edi	f)	
OF	Cancel	Help	

Figure 126. Boîte de dialogue Créer/Editer la tâche d'exploration

11. Cliquez sur le bouton **Nouveau** (dans le coin supérieur droit) pour créer une seconde tâche d'exploration et spécifiez *Down Search* comme nom de tâche dans la boîte de dialogue Nouveaux paramètres.

Create/Edit Mining Task: resp	onse[1]		X
Load Settings: Down Search	•	New 🗙	
🎯 Target Field: 💧	D® response	Target Value: 1	
Simple Settings			
Find segments with:		Low Probability 🔝	
Maximum number of new segments	:	3 🖨	
Minimum segment size			
As percentage of previous segn	ne⊓t (%):	5.0 ≑	
As absolute value (N):		1,000 ≑	
Maximum number of alternatives:		3	
Maximum attributes per segment:		5 🚔	
🗹 Allow attribute re-use within	segment		
Confidence interval for new condition	ons (%):	85.0 🚔	
Expert Settings			
Binning method:	Equal Count	Number of bins:	10
Model search width:	5	Rule search width:	5
Bin merging factor:	2.00		
Allow missing values in conditions:	True	Discard intermediate results:	True
Data			
Build Selection: All Data	- T	1	
Available fields: 🍥 All fields 🔘 🤇	Custom Edi	t	
O	Cancel	Help	

Figure 127. Boîte de dialogue Créer/Editer la tâche d'exploration

- **12**. Faites passer la direction de recherche pour la tâche à **Faible probabilité**. L'algorithme recherchera les segments avec les taux de réponse *les plus faibles* au lieu des plus élevés.
- **13**. Augmentez la taille minimale de segment à 1 000. Cliquez sur **OK** pour revenir au visualiseur Liste interactive.
- 14. Dans le visualiseur Liste interactive, vérifiez que le panneau *Localisateur de segment* affiche les détails de la nouvelle tâche et cliquez sur **Rechercher les segments**.

Segment Finder			
Find segments with:	Low Probability 💌	Settin	gs
Max. no. of new segments:	3 🚔	Find Segments	

Figure 128. Rechercher les segments dans une nouvelle tâche d'exploration

La tâche renvoie un nouvel ensemble d'alternatives, qui est affiché dans l'onglet Alternatives de la boîte de dialogue Albums de modèles et que vous pouvez prévisualiser de la même manière que les résultats précédents.

Name		Target	N	o. of Seamen	ts	Cover	Freq.	Prob.
Alternative 1	1	1		-	3	9,18	3 23:	2 2.539
Alternative 2	1	1			3	9,18	3 233	2 2.539
Alternative 3	1	1			3	8,74	9 14	4 1.659
Alternative Pr	eview							
id	Segmer	nt Rules			Score	Cover (n)	Frequency	Probability
	All segm	ents including	Remainde	r		13,504	1,952	14.45%
1	🖃 mo i m	nths_custon onths_custom	n er er = "0"		1	1,747	0	0.00%
2	⊟ rfm rf	_ score m_score <= 0	.000		1	6,003	0	0.00%
3	⊟ inco in in rf	o me, rfm_sc come > 40297 come <= 5526 m_score > 0.0	ore 1.000 and 17.000 and 100 and		1	1,433	232	16.19%
	rf Remai⊓c	m_score <= 1 ler	0.535			4.321	1.720	39.81%
	Snapshot	13						

Figure 129. Résultats du modèle obtenus par l'intermédiaire de la tâche Down Search

Cette fois, chaque modèle identifie les segments dotés de faibles probabilités de réponse au lieu de fortes probabilités. En examinant la première alternative, vous constatez que le simple fait d'exclure ces segments augmente le taux de correspondance du reste à 39,81%. Ce résultat est inférieur à celui obtenu avec le modèle précédemment étudié, mais il présente une couverture supérieure (et donc un nombre total de correspondances plus élevé).

En combinant les deux approches (utilisation d'une recherche à faible probabilité pour éliminer les enregistrements inintéressants, suivie d'une recherche à forte probabilité), vous pouvez améliorer ce résultat.

15. Cliquez sur **Charger** pour que la première alternative Down Search devienne le modèle de travail et cliquez sur **OK** pour fermer la boîte de dialogue Albums de modèles.

T Je	fake Snapshot at field: Oe response at value: 1			Segment Finder Find segments with Max. no. of new se	gments: 3	Settings
	Segment Rules	Score	Cover (n)		Frequency	Probability
	All segments including Remainder months_customer months_customer = "0"	Excluded		13,504	1,952	14.45%
	rfm_score = 0.000	Excluded		6,003	0	0.00%
	income, rfm_score income > 40297.000 and income <= 55267.000 and rfm_score > 0.000 and rfm_score <= 10.535	1		1,433	232	16.19%
	Remainder			4,321	1,720	39.81%

Figure 130. Exclusion d'un segment

- **16**. Cliquez à droite sur chacun des deux premiers segments et sélectionnez **Exclure le segment**. Ensemble, ces segments capturent presque 8 000 enregistrements avec zéro correspondance entre elles, il est donc souhaitable de les exclure des futures offres. (Les segments exclus auront un score nul pour le signaler.)
- 17. Cliquez avec le bouton droit de la souris sur le troisième segment et sélectionnez **Supprimer le segment**. Le taux de correspondance de 16,19 % de ce segment n'est pas très différent du taux de référence de 14,45 %, et il n'ajoute donc pas assez d'informations pour justifier sa conservation.

Remarque : La suppression d'un segment et son exclusion sont deux opérations différentes. L'exclusion d'un segment modifie uniquement son score, alors que sa suppression le retire complètement du modèle.

Une fois que vous avez exclu les segments ayant les plus basses performances, vous pouvez rechercher les segments avec les plus hautes performances dans le reste.

18. Dans la table, cliquez sur la ligne du reste pour la sélectionner de telle manière que la prochaine tâche d'exploration s'applique uniquement au reste.

nte File	ractive List: response[1] #2 ≧ Edit	8 🕒 🐚 😭	(*) (#) (*) (*)]
if 1 arge arge	Take Snapshot et field: O response et value: 1		Segment Finder Find segments with Max. no. of new se	h: Low Probability egments: 3	Settings
	Segment Rules	Score	Cover (n)	Frequency	Probability
	All segments including Remainder		13,504	4 1,952	14.45%
P	months_customer months_customer = "0"	Excluded	1,747	0	0.00%
2		Excluded	6,003	3 0	0.00%
	Remainder		5,754	4 1,952	33.92%

Figure 131. Sélection d'un segment

- **19**. Le reste étant sélectionné, cliquez sur **Paramètres** pour ouvrir à nouveau la boîte de dialogue Créer/Editer une tâche d'exploration.
- 20. En haut, dans Charger les paramètres, sélectionnez la tâche d'exploration par défaut : response[1].
- **21**. Modifiez les **Paramètres simples** pour augmenter le nombre de nouveaux segments jusqu'à 5 et la taille minimale de segments à 500.
- 22. Cliquez sur OK pour revenir au visualiseur Liste interactive.

Create/Edit Mining Task: Dow	n Search		X
Load Settings: response[1]		New 🗙	
🎯 Target Field:	🔎 response	Target Value: 1	
Simple Settings			
Find segments with:		High Probability 🔻	
Maximum number of new segments	5:	5	
Minimum segment size			
As percentage of previous segr	ment (%):	5.0 ≑	
As absolute value (N):		500 🚔	
Maximum number of alternatives:		3	
Maximum attributes per segment:		5	
📝 Allow attribute re-use within	i segment		
Confidence interval for new conditi	ions (%):	85.0	
Expert Settings			
Binning method:	Equal Count	Number of bins:	10
Model search width:	5	Rule search width:	5
Bin merging factor:	2.00		
Allow missing values in conditions:	True	Discard intermediate results:	True
Data			
Build Selection: All Data	T 🖂	1	
Available fields: 🔘 All fields 🔘 0	Custom Edi	t)	
O	Cancel	Help	

Figure 132. Sélection de la tâche d'exploration par défaut

23. Cliquez sur Rechercher les segments.

Cette action affiche encore un nouvel ensemble de modèles alternatifs. En insérant les résultats d'une tâche d'exploration dans une autre tâche, ces derniers modèles contiennent un mélange de segments très performants et de segments peu performants. Les segments dotés de taux de réponse faibles sont exclus, ce qui signifie que leur score est nul, alors que les segments inclus ont le score 1. Les statistiques globales reflètent ces exclusions, avec un taux de correspondance de 45,63 % pour le premier modèle alternatif et une couverture supérieure (1 577 correspondances sur 3 456 enregistrements) à celle de tous les modèles précédents.

Name	Target	No. of Segme	nts	Cover			Freq.	Prob.	
Alternative	1 1		7		3,	456	1,5	577	45.63
Alternative	2 1		7	-	3,	456	1,5	577	45.63
Alternative	3 1		7		3,	456	1,8	577	45.63
Alternative	e Preview								
id	Segment Rules		Score	(Cover (n)	Free	quency	Probability	(
	All segments including Remaind	er			13,504		1,952	14.4	5%
1	months_customer months_customer = "0"		Excluded	1	1,747		0	0.0	0%
2	⊡ rfm_score rfm_score <= 0.000		Excluded	ı	6,003	1	0	0.0	0%
3	☐ rfm_score, income rfm_score > 12.333 and income > 52213.000		1		555	i	456	82.1	6%
4	□ income income > 55267.000		1		643	,	551	85.6	9%
5	□ number_transactions, rf number_transactions > 2 rfm_score > 12.333	m_score .000 and	1		533	1	206	38.6	5%
	Snapshots								

Figure 133. Alternatives pour un modèle associé

24. Prévisualisez la première alternative et cliquez sur Charger pour en faire le modèle de travail.

Calcul des mesures personnalisées avec Excel

1. Pour avoir une meilleure visibilité sur la façon dont le modèle fonctionne en termes pratiques, choisissez **Organiser les mesures du modèle** dans la barre d'outils.

Inte	eractive List: respons	e[1] #4					
🔰 <u>F</u> ile	e 📄 <u>E</u> dit 💰 <u>V</u> iew	Tools 🖏 Generate 🛛 🗐 🛃			d) 📭 🐓		0
/iewe	r Gains Annotations	Find Segments Settings					
	Take Snapshot Organize Model Measures				Segment Finder		
Targ	et field: 🔘 response	Image: Organize Data Selections Image: Organize Data Selections Image: Organize Data Selections	~		Find segments wr Max. no. of new s	th: High Probability segments: 5	Settings
Targ	et value: 1	💕 Take Snapshot					
id	Segment Rules		Score	Cover (n)	Frequency	Probability
	All segments including Re	mainder			13,504	1,952	14.45%
1	months_customer months_customer	= "0"	Excluded		1,747	0	0.00%
2	☐ rfm_score rfm_score <= 0.000	D	Excluded		6,003	0	0.00%
3	■ rfm_score, income rfm_score > 12.33 income > 52213.00	; 3 and 0	1		555	456	82.16%
4	income > 55267.00	0	1		643	551	85.69%
5	number_transaction number_transaction rfm_score > 12.333	e ns, rfm_score ns > 2.000 and 3	1		533	206	38.65%
Mode	J Summary; Cover 3,456: Fi	requency 1,577: Probability 45.63%					<u> </u>

Figure 134. Organisation des mesures de modèle

La boîte de dialogue Organiser les mesures du modèle vous permet de choisir les mesures (ou colonnes) à afficher dans le visualiseur Liste interactive. Vous pouvez aussi indiquer si les mesures sont calculées sur tous les enregistrements ou sur un sous-ensemble sélectionné, et vous pouvez choisir d'afficher un graphique circulaire plutôt qu'un nombre le cas échéant.

over Coverage Pie Chart All Data over (n) Coverage Numeric All Data requency Frequency Numeric All Data orobability Probability Numeric All Data over Frequency Numeric All Data orobability Probability Numeric All Data over Error Numeric All Data Custom Measures Calculate custom measures in Excel (TM): Yes No Connect to Excel (TM) Workbook: Name Description	lame	Туре	Display	Data Selection	Show
over (n) Coverage Numeric All Data requency Frequency Numeric All Data robability Probability Numeric All Data rror Error Numeric All Data Custom Measures Calculate custom measures in Excel (TM): Image Yes Image No Connect to Excel (TM) Workbook:	over	Coverage	Pie Chart	All Data	
equency Frequency Numeric All Data I and All Data I and All Data I and All Data I and I an	over (n)	Coverage	Numeric	All Data	
Probability Probability Numeric All Data ror Error Numeric All Data Custom Measures Calculate custom measures in Excel (TM): Image: Second	requency	Frequency	Numeric	All Data	
ror Error Numeric All Data	robability	Probability	Numeric	All Data	
Custom Measures Calculate custom measures in Excel (TM): Yes No Connect to Excel (TM) Workbook: Name Description Show	rror	Error	Numeric	All Data	
	Custom Measures Calculate custom mea Connect to Excel (TM	isures in Excel (TM): @ Ye	s 🔘 No		
	Custom Measures Calculate custom mea Connect to Excel (TM Name	isures in Excel (TM): Ye Munu Workbook: Description	s 🔘 No		Show
	Custom Measures Calculate custom mea Connect to Excel (TM Name	isures in Excel (TM): Ye Workbook: Description	s 🔘 No		Show

Figure 135. Boîte de dialogue Organiser les mesures du modèle

En outre, si Microsoft Excel est installé, vous pouvez lier un modèle Excel qui calculera les mesures personnalisées et les ajoutera à l'affichage interactif.

- 2. Dans la boîte de dialogue Organiser les mesures du modèle, configurez **Calculer les mesures** personnalisées dans Excel (TM) sur Oui.
- 3. Cliquez sur Se connecter à MS Excel (TM)
- 4. Sélectionnez le classeur *template_profit.xlt*, situé dans le répertoire des *flux* dans le dossier *Demos* de votre installation IBM SPSS Modeler et cliquez sur **Ouvrir** pour lancer la feuille de calcul.

	ticro	soft Ex	cel - templat	e_profit1							1 💌
:2	Eile	<u>E</u> dit	<u>V</u> iew <u>I</u> nsert	F <u>o</u> rmat <u>T</u> oo	s <u>D</u> ata	<u>W</u> indow <u>H</u>	lelp Adol	<u>b</u> e PDF		- 8	×
	₽ ₽	Arial		• 10 • B	ΙU		·•• 🛒	% ,	🚛 🖽 🗸	🗞 • <u>A</u> •	++ ∓
: 🔁	1	-									
	F4		▼ fx	=IF(H4="",0,L	4)-Settin	gs!FIX_1					
100	A	В	C	D		E			F	G	^
											=
1											
2											
			Metric:	Imported	Metric:	Calculate	d Metric:	Calculate	d Metric:		
3	#	Use	Frequency	Cover		Profit Mar	gin	Cumulativ	e Profit	larget	-
4	1								-2,500.00		
										-	
5	2										
14 4	+ 1		lel Measures	/ Settings /	Configura	ation /	< .			1	>
1400000											

Figure 136. Feuille de calcul Excel de modèles de mesures

Le modèle Excel contient trois feuilles de calcul :

- Mesures du modèle affiche les mesures du modèle importées du modèle et calcule les mesures personnalisées pour les réexporter vers le modèle.
- Paramètres contient les paramètres à utiliser dans le calcul des mesures personnalisées.
- Configuration définit les mesures à importer du modèle et à exporter vers ce modèle.

Les mesures réexportées vers le modèle sont :

- Marge de profit. Revenus nets du segment
- Profit cumulé. Total des profits de la campagne

Définis par les formules suivantes :

Marge de profit = Fréquence * Revenu par personne sondée - Couverture * Coût variable

Profit cumulé = Marge de profit totale - Coût fixe

Notez que la fréquence et la couverture sont importées du modèle.

Les paramètres de coût et de revenus sont indiqués par l'utilisateur dans la feuille de calcul Paramètres.

Bile Edit Yew Insert Figmat Lools Data Window Help Adobe PDF - 8 × Arial 10 B I	X	licro	soft E	xcel -	templa	te_p	rofit1	1											-	
Arial 10 B U E E W Settings % Settings % Settings Configuration K W K <td>:2</td> <td>Eile</td> <td>Edit</td> <td><u>V</u>iew</td> <td>Insert</td> <td>Fo</td> <td>rmat</td> <td>Too</td> <td>ls</td> <td><u>D</u>ata</td> <td>Window</td> <td>v <u>H</u>elp</td> <td>Ade</td> <td>o<u>b</u>e PDF</td> <td></td> <td></td> <td></td> <td></td> <td>_ é</td> <td>7 X</td>	:2	Eile	Edit	<u>V</u> iew	Insert	Fo	rmat	Too	ls	<u>D</u> ata	Window	v <u>H</u> elp	Ade	o <u>b</u> e PDF					_ é	7 X
J12 fx A B D E F G H I 4 -	:0	12 I	Arial			v 10	•	в	I	U		=	9	%	, <u>*.</u> 0	*		· 37 -	A	* 12
J12 fx A B D E F G H I 4	-	1	*]_																	_
A B D E F G H I 4		J12	-	-	f _×															
4				A			E	В		D		Е		F	G		Н		1	~
5	4										_				_					_
6	5										12		_		-					_
7	6					-			_				-		_					- (***
8	7								_		9		_		-					-11
9 Image: strain of the st	8					_			_									-		-11
10 11 12 Costs and revenue 12 Costs and revenue 13 - 14 - 14 - 14 - 14 - 15 - 100.00 14 - 100.00 14 - 100.00 15 - 100.00 16 17 17 17 17 16 16 16 17 <	9								-											- 1
12 Costs and revenue	10					-			-				-		-					- 11
12 Costs and revenue 2,500.00 Image: Cost of the second se	11	Cart	and						-		15				2					-11
13 - Tried costs 2,00.00 14 - Variable cost 0.50 15 - Revenue per respondent 100.00 16 - 17 - 18 - 19 - 20 - 21 - IM Model Measures \ Settings \ Configuration \	12	Eix	s and od.cov	reven	ue		2	500	nn											= =
17 0.50 15 - Revenue per respondent 16 100.00 17 100.00 18 100.00 19 100.00 20 100.00 21 100.00 14 → H\ Model Measures \ Settings \ Configuration /	14	- Var	iahle i	rnet		-1	۷,	.000. N	50		12		-							-
16 17 17 18 19 20 21 ▶ N Model Measures Settings Configuration /	15	- Rev	enue	ner res	soonder	nt		100	nn											- 11
17 18 19 20 21 ✓ 14 ← ▶ N Model Measures Settings Configuration / <	16	1.0.	0.1.00												1					
18 19 20 21 IM IM IM IM Settings (Configuration / K mm)	17																			
19 20 21 IM IM <	18																			
20 21 Model Measures Settings Configuration / ✓	19																			
21 I	20																			
	21								C = 1	6	-	r							-	
5 U 15 A	Dec.	••	M / M	oael M	easures	λ 26	etting	JS (Cor	mgur	ation /	1	•						13	

Figure 137. Feuille de calcul Excel Paramètres

Coût fixe est le coût configuré pour la campagne ; par exemple, conception et planification. **Coût variable** est le coût d'extension de l'offre à chaque client, par exemple les enveloppes et les timbres.

Recettes par personne sondée est le revenu net d'un client qui répond à l'offre.

5. Pour terminer la liaison de retour vers le modèle, utilisez la barre des tâches Windows (ou appuyez sur Alt+Tab) pour revenir au visualiseur Liste interactive.

Hint: Use this dialog to cho nputs to calculate custom	ose which model measures will be used by Excel (TM) as measures.
Input	Model Measure
Frequency	Frequency
Cover	Cover (n)

Figure 138. Choix des entrées de mesures personnalisées

La boîte de dialogue Choisir les entrées de mesures personnalisées s'affiche, vous permettant de faire correspondre les entrées du modèle aux paramètres spécifiques définis dans le modèle. La colonne de gauche répertorie les mesures disponibles et la colonne de droite les fait correspondre aux paramètres de la feuille de calcul définis dans la feuille de calcul Configuration.

6. Dans la colonne **Mesures de modèle**, sélectionnez **Fréquence** et **Couverture (n)** pour les entrées respectives puis cliquez sur OK.

Dans ce cas, les noms du paramètre du modèle (Fréquence et Couverture (n)) correspondent aux entrées, mais des noms différents peuvent aussi être utilisés.

7. Cliquez sur **OK** dans la boîte de dialogue Organiser les mesures du modèle pour mettre à jour le visualiseur Liste interactive.

over	10 X 1920	Display	Data Selection	Show
519 T 011	Coverage	Pie Chart	All Data	
over (n)	Coverage	Numeric	All Data	
requency	Frequency	Numeric	All Data	
robability	Probability	Numeric	All Data	
rror	Error	Numeric	All Data	
Connect to Excel (Th	asures in Excel (TM): () Ye ()) Workbook: Files\SF	s (O) No PSSInc\PASVVMode	ler14\Demos\Classification_Mod	uleitemplate_profit.xtt
Connect to Excel (Th	asures in Excel (TM): O Ye (M) Workbook: Files/SF	s () No PSSInc/PASV/Mode	ler14\Demos\Classification_Mod	ulettemplate_profit.xtt
Connect to Excel (Th Name Profit margin	asures in Excel (TM):	s () No PSSInc'PASV/Mode	ler14\Demos\Classification_Mod	ulettemplate_profit.xtt Show

Figure 139. Boîte de dialogue Organiser les mesures du modèle avec les mesures personnalisées d'Excel

Les nouvelles mesures sont maintenant ajoutées en tant que nouvelles colonnes dans la fenêtre et seront recalculées chaque fois que le modèle sera mis à jour.

wer Irge Irge	Gains Annotations ake Snapshot It field: Or response It value: 1		-Segme Find se Max. no	nt Finder gments with: High Pr b. of new segments:	robability 🔻	Find S	Settings
	Segment Rules	Score	Cover (n)	Frequency	Probability	Profit margin	Cumulative
	All segments including Remainder		13,504	1,952	14.45%	0	0 🗲
	months_customer months_customer = "0"	Excluded	1,747	0	0.00%	-873.5	-2,500
2	□ rfm_score rfm_score <= 0.000	Excluded	6,003	0	0.00%	-3,001.5	-2,500
3	☐ rfm_score, income rfm_score > 12.333 and income > 52213.000	1	555	456	82.16%	45,322.5	42,822.5
ł	☐ income income > 55267.000	1	643	551	85.69%	54,778.5	97,601
5	number_transactions, rfm_score number_transactions > 2,000 and rfm_score > 12,333	1	533	206	38.65%	20,333.5	117,934.5
del	Summary; Cover 3,456: Frequency 1,577: Probabil	ity 45.63%					V

Figure 140. Mesures personnalisées d'Excel affichées dans le visualiseur Liste interactive

En éditant le modèle Excel, vous pouvez créer autant de mesures personnalisées que vous le souhaitez.

Modification du modèle Excel

Bien qu'IBM SPSS Modeler propose un modèle Excel par défaut à utiliser avec le visualiseur Liste de décisions, il est possible de modifier les paramètres ou d'ajouter les vôtres. Par exemple, les coûts dans le modèle peuvent ne pas correspondre à ceux de votre entreprise et doivent être modifiés.

Remarque : Si vous modifiez un modèle existant, ou que vous créez le vôtre, n'oubliez pas d'enregistrer le fichier avec un suffixe *.xlt* d'Excel 2003.

Pour modifier le modèle par défaut et y ajouter de nouveaux coûts et informations sur les revenus et mettre à jour le visualiseur Liste interactive en y ajoutant de nouveaux chiffres :

- 1. Dans le visualiseur Liste interactive, sélectionnez **Organiser les mesures de modèle** dans le menu Outils.
- 2. Dans la boîte de dialogue Organiser les mesures du modèle, cliquez sur **Connecter à Excel**[™].
- 3. Sélectionnez le classeur *template_profit.xlt* et cliquez sur **Ouvrir** pour lancer la feuille de calcul.
- 4. Sélectionnez la feuille de calcul Paramètres.
- 5. Modifiez les coûts fixes sur 3250,00 et le Revenu par personne interrogée sur 150,00.

	licros	oft E	kcel -	templa	te_pr	ofit1.xlt							6	-)0	×
:8	Eile	Edit	⊻iew	Insert	For	nat <u>T</u> ools	Data	<u>W</u> indow	Help	Ado <u>b</u> e PDF				- 8	x
1	21	Arial			v 10	- B I	U		-a-	🥶 % ,	*.0	=	ð -	A -	
-	-														
	GAIN	1	-	fx	150										
		-	A			В	D	E		F	G	H			~
4															
5															
6															
7															
8															
9															-
10															-
11	8/072 - 30	0.9						11							-
12	Cost	s and	reven	nue											- =
13	- Fixe	ed cos	sts			3,250.00	-								_
14	- Var	iable (cost			0.50	-			_					-
15	- Rev	enue	per re:	sponder	nt	150.00									-
16					-										-
17															-
18															-
19															-
20															~
H A		I M	odel M	leasures	Set	ttings / Co	nfigur	ation /	<		1111			>	
Read	ly											NUM			

Figure 141. Valeurs modifiées sur la feuille de calcul Excel Paramètres

6. Sauvegardez le modèle modifié en utilisant un nom de fichier unique et approprié. Vérifiez qu'il possède une extension *.xlt* d'Excel 2003.

Save As			? 🛛
Save in:	🛅 Classifica	tion_Module 🛛 🕑 🕶 🖄 🔕 🗙 📸 🎹 🕶 To	ooļs 🕶
My Recent Documents	template_p template_p	rofit1.xlt rofit.xlt	
My Documents			
My Computer			
Mu Network	File <u>n</u> ame:	template_profit_3250l,xlt	Save
Places	Save as <u>t</u> ype:	Template (*.xlt)	Cancel

Figure 142. Enregistrement d'un modèle Excel modifié

7. Utilisez la barre des tâches de Windows (ou appuyez sur Alt+Tab) pour retourner au visualiseur Liste interactive.

Dans la boîte de dialogue Choisir les entrées de mesures personnalisées, sélectionnez les mesures à afficher et cliquez sur **OK**.

8. Cliquez sur **OK** dans la boîte de dialogue Organiser les mesures du modèle pour mettre à jour le visualiseur Liste interactive.

Bien sûr, cet exemple ne présente qu'une seule façon de modifier le modèle Excel. Vous pouvez effectuer d'autres modifications qui extraient des données du visualiseur Liste Interactive ou qui lui transmettent des données, ou travailler depuis Excel pour produire d'autres entrées, tels que des graphiques.

n arg arg	Take Snapshot et field: O response et value: 1		- Seg Find Max	ment Finder segments with: High . no. of new segments:	n Probability 🤝	Find	Segments
d	Segment Rules	Score	Cover (n)	Frequency	Probability	Profit margin	Cumulative
	All segments including Remainder		13,504	1,952	14.45%	0	0 4
1	months_customer months_customer = "0"	Excluded	1,74;	0	0.00%	-873.5	-3,250
2	l ⊟ rfm_score rfm_score ≺= 0.000	Excluded	6,003	3 0	0.00%	-3,001.5	-3,250
3	☐ rfm_score, income rfm_score > 12.333 and income > 52213.000	1	55	5 456	82.16%	68,122.5	64,872.5
4	☐ income income > 55267.000	1	643	3 551	85.69%	82,328.5	147,201
5	number_transactions, rfm_score number_transactions > 2.000 and rfm_score > 12.333	1	53:	3 206	38.65%	30,633.5	177,834.5

Figure 143. Mesures personnalisées d'Excel modifiées affichées dans le visualiseur Liste interactive

Enregistrement des résultats

Pour enregistrer un modèle afin de pouvoir l'utiliser ultérieurement au cours de la session interactive, vous pouvez prendre un instantané du modèle, qui apparaîtra dans l'onglet Instantanés. Vous pouvez accéder aux instantanés enregistrés à tout moment au cours de la session interactive.

Ainsi, vous pouvez tester d'autres tâches d'exploration pour rechercher des segments supplémentaires. Vous pouvez également éditer des segments existants, insérer des segments personnalisés sur la base de vos propres règles métier, créer des sélections de données pour optimiser le modèle pour des groupes précis et personnaliser le modèle de différentes manières. Enfin, vous pouvez inclure ou exclure explicitement chaque segment, selon vos besoins, pour préciser comment chacun d'eux sera évalué.

Lorsque vous êtes satisfait des résultats, utilisez le menu Générer pour générer un modèle qui peut être ajouté aux flux ou déployé à des fins d'évaluation.

Une autre solution pour enregistrer l'état actuel de la session interactive et y revenir un autre jour consiste à choisir **Mettre à jour le noeud de modélisation** dans le menu Fichier. Ainsi, le noeud de modélisation Liste de décision sera mis à jour avec les paramètres en cours, y compris les tâches

d'exploration, les instantanés de modèle, les sélections de données et les mesures personnalisées. Lorsque vous réexécutez le flux, vérifiez que l'option **Utiliser les informations de session interactive enregistrées** est sélectionnée dans le noeud de modélisation Liste de décision pour restaurer l'état actuel de la session.

Chapitre 12. Classification des clients de télécommunications (régression logistique multinomiale)

La régression logistique est une technique statistique de classification des enregistrements sur la base des valeurs des champs d'entrée. Excepté le fait qu'elle utilise un champ cible catégoriel et non pas numérique, cette régression est semblable à la régression linéaire.

Par exemple, supposons qu'un fournisseur de télécommunications ait segmenté sa base de clientèle par modèles d'utilisation de service, classant ses clients en quatre groupes. Si les données démographiques peuvent être utilisées pour prévoir les affectations de groupes, vous pouvez personnaliser les offres pour les clients éventuels.

Cet exemple utilise le flux *telco_custcat.str*, qui fait référence au fichier de données *telco.sav*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *telco_custcat.str* se trouve dans le répertoire des *flux*.

Cet exemple est axé sur l'utilisation des données démographiques dans le but de prévoir des modèles d'utilisation. Le champ cible *custcat* possède quatre valeurs possibles qui correspondent aux quatre groupes de clients suivants :

Valeur	Libellé
1	Basic Service
2	E-Service
3	Plus Service
4	Total Service

Comme le champ cible contient plusieurs catégories, un modèle multinomial est utilisé. Dans le cas d'un champ cible comprenant deux catégories distinctes, telles que oui/non, vrai/faux ou attrition/absence d'attrition, un modèle binomial peut être créé. Pour plus d'informations, voir la rubrique Chapitre 13, «Attrition dans le domaine des télécommunications (régression logistique binomiale)», à la page 141.

Création du flux

1. Ajoutez un noeud source de fichier Statistics pointant vers telco.sav dans le dossier Demos.



Figure 144. Flux d'échantillons permettant de classifier les clients par régression logistique multinomiale

a. Ajoutez un noeud type et cliquez sur **Lire les valeurs**, en vous assurant que tous les niveaux de mesure sont correctement paramétrés. Par exemple, la majorité des champs avec des valeurs 0 et 1 peuvent être considérés comme des champs indicateurs.

😡 Туре								×
Types Format	Annotations						0	-
~	PRE	ad Valu	es Clea	r Values	CI	lear All Val	ues	
Field -	Measurem	ent	Values	Missing		Check	Ro	le
🔆 gender	💑 Nominal		0,1		No	ne	🔪 Inpu	.t 🖆
🔆 reside	🖉 Continuous		[1,8]		No	ne	🔪 Inpu	.t
🗘 tollfree	🖁 Flag		1/0		No	ne	🔪 Inpu	.t
📿 equip	🖁 Flag		1./0		No	ne	🖒 Inpu	.t
📿 callcard	🎖 Flag		1/0		N	Defeut		
📿 wireless	🔓 Flag 🛛 🗖		1/0		NIC	Sperault	ć .	
🛞 longmon	🔗 Continuou	Se	ect All			Continuo	us	
🛞 tollmon	Scontinuou	Se	ect None			Categori	cal	
View current	fields OVi	Se	ect Fields			Flag	Þ	
OK Cance		Co	py ste Special	Ctrl+C . Ctrl+V		Nominal Ordinal	-10	Reset

Figure 145. Configuration du niveau de mesure pour plusieurs champs

Conseil : Pour modifier les propriétés de plusieurs champs contenant des valeurs similaires (telles que 0/1), cliquez sur l'en-tête de colonne *Valeurs* afin de trier les champs en fonction de cette valeur. Maintenez la touche Maj enfoncée tout en utilisant la souris ou les touches fléchées pour sélectionner tous les champs à modifier. Cliquez ensuite sur la sélection avec le bouton droit de la souris pour modifier le niveau de mesure ou les autres attributs des champs sélectionnés.

Veuillez noter que puisqu'il est plus correct de considérer le *sexe* comme un champ avec un ensemble de deux valeurs plutôt que comme un indicateur, laissez sa valeur de mesure sur **Nominal**.

b. Définissez le rôle du champ *custcat* sur **Cible**. Le rôle de tous les autres champs doit être défini sur **Entrée**.

Type	Preview				0
Types Format	Annotations	ilues Clear	Values	Clear All Va	alues
Field -	Measurement	Values	Missing	Check	Role
epili	🕘 гіаў	170	-	NONE	a input
🐡 loglong	Continuous	[-0.10536		None	> Input
🔅 logtoll	Continuous	[1.74919		None	🔪 Input
🛞 logequi	🔗 Continuous	[2.73436		None	🔪 Input
🋞 logcard	🔗 Continuous	[1.01160		None	🔪 Input
🛞 logwire	🔗 Continuous	[2.70136		None	🔪 Input
🛞 Ininc	🔗 Continuous	[2.19722		None	🔪 Input
🔆 custcat	💑 Nominal	1,2,3,4		None	O Target
A also we	💑 Nominal	0,1		None	🔪 Input

Figure 146. Définition du rôle de champ

Cet exemple étant axé sur les données démographiques, utilisez un noeud Filtrer pour n'inclure que les champs pertinents (*region, age, marital, address, income, ed, employ, retire, gender, reside* et *custcat*). Les autres champs peuvent être exclus pour cette analyse.

Demographic Preview Fitter Annotations		
7.	Fields:	42 in, 31 filtered, 0 renamed, 11 ou
Field -	Filter	Field
region	\rightarrow	region
tenure	× >	tenure
age	\rightarrow	age
marital	\rightarrow	marital
address	\rightarrow	address
income	\rightarrow	income
ed	\rightarrow	ed
employ	\rightarrow	employ
retire	\rightarrow	retire
gender	\rightarrow	gender 📃 🔽
View current fields View OK Cancel	w unused field	settings

Figure 147. Filtrage des champs démographiques

(Vous pouvez également paramétrer le rôle sur **Aucun** pour ces champs plutôt que de les exclure, ou sélectionner les champs que vous souhaitez utiliser dans le noeud de modélisation.)

2. Dans le noeud Logistique, cliquez sur l'onglet **Modèle** et sélectionnez la méthode **Pas à pas**. Sélectionnez **Multinomial**, **Effets principaux** et **Inclure la constante dans l'équation**.

😡 custcat			X
		0 -	
Fields Model Expert Analy	ze Annotations		
Model name: 💿 Auto 🛇 Custi	om		
👿 Use partitioned data			
👿 Build model for each split			
Procedure: 💿 Multinomial	C) Binomial	
Multinomial Procedure			-
Method: St	epwise		
Base category for target: 1	Specify.		
Model type: 🔘 Main Effects	O Full Factorial	🔘 Custom	
Model Terms:			
		×	
Vinclude constant in equation			
OK 🕨 Run Cancel)	Apply Re:	set

Figure 148. Choix des options de modèle

Laissez la catégorie de base de la cible définie sur 1. Le modèle comparera les autres clients à ceux qui sont abonnés au Basic Service.

3. Dans l'onglet Expert, sélectionnez le mode **Expert**, puis **Sortie** et, dans la boîte de dialogue Sorties avancées, sélectionnez **Matrice de confusion**.

💟 Logistic Regression: Advanced Output 🛛 🛛 🔀					
Summary statistics	🔲 Parameter estimates				
🔲 Likelihood ratio tests	Confidence interval:	95.0 ≑			
Asymptotic correlation	📃 Asymptotic covariance				
🔲 Goodness of fit chi-square statistics	🗹 Classification table				
lteration history for every	1 🖨	step(s)			
📃 Stepwise variable loadings	📕 Monotonicity measures				
Information criteria					
OK Cancel Help					

Figure 149. Choix des options de sortie

Navigation dans le modèle

1. Exécutez le noeud pour générer le modèle, qui est ajouté à la palette Modèles dans l'angle supérieur droit. Pour afficher ses détails, cliquez avec le bouton droit de la souris sur le noeud du modèle généré et sélectionnez **Parcourir**.
L'onglet Modèle affiche les équations utilisées pour affecter les enregistrements à chaque catégorie du champ cible. Il existe quatre catégories possibles, l'une d'elles est la catégorie de base pour laquelle aucun détail d'équation ne s'affiche. Les détails sont affichés pour les trois équations restantes, où la catégorie 3 représente le Plus Service et ainsi de suite.



Figure 150. Navigation dans les résultats du modèle

L'onglet Récapitulatif affiche (entre autres) la cible et les entrées (champs prédicteurs) utilisées par le modèle. Ces champs sont ceux qui ont été réellement choisis sur la base de la méthode Pas à pas, et non la liste complète soumise.



Figure 151. Récapitulatif du modèle avec champs cible et champs d'entrée

Les éléments affichés dans l'onglet Options avancées dépendent des options sélectionnées dans la boîte de dialogue Sorties avancées, dans le noeud de modélisation.

L'élément Récapitulatif du traitement des observations est systématiquement affiché. Il indique le pourcentage d'enregistrements inclus dans chaque catégorie du champ cible. Vous pouvez ainsi utiliser un modèle nul servant de base à la comparaison.

Sans créer de modèle qui utilise des prédicteurs, votre meilleure prévision consiste à affecter tous les clients au groupe le plus commun, le groupe du Plus Service.

custcat	t	Previ	ew 📳 🗿 🕞	. (r
Aodel Su	immary Advanced Settings Annotat	ions		
۲				
ase Proc	essing Summary			
ase Proc	essing Summary	N	Marginal Percentage	
ase Proc	essing Summary Basic service	N 266	Marginal Percentage 26.6%	
ase Proc	essing Summary Basic service E-service	N 266 217	Marginal Percentage 26.6% 21.7%	
ase Proc custcat	essing Summary Basic service E-service Plus service	N 266 217 281	Marginal Percentage 26.6% 21.7% 28.1%	
ase Proc	essing Summary Basic service E-service Plus service Total service	N 266 217 281 236	Marginal Percentage 26.6% 21.7% 28.1% 23.6%	

Figure 152. Récapitulatif du traitement des observations

En fonction des données d'apprentissage, si vous avez affecté tous les clients au modèle nul, votre prévision est correcte 281/1000 = 28,1 % du temps. L'onglet Options avancées contient des informations supplémentaires qui vous permettent d'examiner les prévisions du modèle. Vous pouvez ensuite comparer les prévisions aux résultats du modèle nul pour voir le fonctionnement de votre modèle avec vos données.

En bas de l'onglet Options avancées, la table de classification supervisée affiche les résultats de votre modèle, qui est correct 39,9 % du temps.

Votre modèle est particulièrement performant à l'heure d'identifier les clients Total Service (catégorie 4), mais fonctionne très mal pour l'identification des clients E-service (catégorie 2). Si vous souhaitez une meilleure exactitude pour les clients de la catégorie 2, vous devez trouver un autre prédicteur pour les identifier.

🖓 custcat 🛛 🔀										
Model Summary A	Model Summary Advanced Settings Annotations									
	Predicted									
Observed	Basic Plus Total Percent Observed service E-service service service Correct									
Basic service	122	8	75	61	45.9%					
E-service	E-service 58 10 68 81 4.6%									
Plus service	89	8	133	51	47.3%					
Total service	47	12	43	134	56.8%					
Overall Percentage	31.6%	3.8%	31.9%	32.7%	39.9%					
OK Cancel	OK Cancel									

Figure 153. Table de classification

En fonction de ce que vous souhaitez prévoir, le modèle peut s'avérer parfaitement adapté à vos besoins. Par exemple, si l'identification des clients de la catégorie 2 ne vous intéresse pas, le modèle peut être assez précis pour vous. Cela peut être le cas lorsque le E-service est un produit d'appel qui ne génère que peu de bénéfices.

Si, par exemple, votre plus grand retour sur investissement provient des clients des catégories 3 ou 4, il est possible que le modèle vous fournisse les informations nécessaires.

Pour évaluer le niveau d'adéquation du modèle aux données, divers diagnostics sont disponibles dans la boîte de dialogue Sorties avancées lorsque vous créez le modèle. Des explications sur les fondements mathématiques des méthodes de modélisation utilisées dans IBM SPSS Modeler sont présentées dans le *guide des algorithmes d'IBM SPSS Modeler*, disponible dans le répertoire *Documentation* du disque d'installation.

Sachez également que ces résultats sont établis uniquement d'après les données d'apprentissage. Pour évaluer comment le modèle peut se généraliser à d'autres données dans le monde réel, vous pouvez utiliser un noeud Partitionner destiné à contenir un sous-ensemble d'enregistrements à des fins de test et de validation.

Chapitre 13. Attrition dans le domaine des télécommunications (régression logistique binomiale)

La régression logistique est une technique statistique de classification des enregistrements sur la base des valeurs des champs d'entrée. Excepté le fait qu'elle utilise un champ cible catégoriel et non pas numérique, cette régression est semblable à la régression linéaire.

Cet exemple utilise le flux *telco_churn.str*, qui fait référence au fichier de données *telco.sav*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *telco_churn.str* se trouve dans le répertoire des *flux*.

Par exemple, supposons qu'un fournisseur de télécommunications souhaite connaître le nombre de clients qui partent à la concurrence. Si les données d'utilisation du service permettent de prédire les clients susceptibles de passer à un autre fournisseur, les offres peuvent être personnalisées afin de retenir autant de clients que possible.

Cet exemple explique comment se servir des données d'utilisation pour prédire la perte de clients (attrition). Etant donné que la cible présente deux catégories distinctes, un modèle binomial est utilisé. Si la cible présente plus de deux catégories, un modèle multinomial peut être créé à la place. Pour plus d'informations, voir la rubrique Chapitre 12, «Classification des clients de télécommunications (régression logistique multinomiale)», à la page 133.

Création du flux

1. Ajoutez un noeud source de fichier Statistics pointant vers telco.sav dans le dossier Demos.



Figure 154. Flux d'échantillons permettant de classifier les clients par régression logistique binomiale

2. Ajoutez un noeud type pour définir des champs, en vous assurant que tous les niveaux de mesure sont correctement paramétrés. Par exemple, la plupart des champs dont les valeurs sont 0 et 1 peuvent être considérés comme des champs indicateurs. Cependant, certains champs, tels que celui indiquant le genre, doivent être considérés comme des champs nominaux à deux valeurs.

Types Format	Annotations						0	
×- 00 0	P Re	ad Value	es Clear	Values	Cle	ear All Val	ues	
Field -	Measurem	ent	Values	Missing		Check	F	Role
🔆 gender	💑 Nominal		0,1		Nor	ne		put 📥
🚫 reside	🔗 Continuous		[1,8]		Nor	ne	N In	put
🗘 tollfree	🎖 Flag		1/0		Nor	ne	🔪 Inj	put 📑
📿 equip	🎖 Flag		1/0		Nor	ne	🔪 Inj	put 🚽
📿 callcard	🎖 Flag		1/0		No	-Defeult		
📿 wireless	🖉 Flag 🗖		1/0		NI	Speradic	5	
🛞 longmon	🔗 Continuou	Sel	ect All			Continuo	us	
🐲 tolimon	🔗 Continuou	Sel	ect None			Categori	cal	
View current	fields OVi	Sel	ect Fields	•	•	Flag	D	
OK Cancel		Col	py ste Special	Ctrl+C . Ctrl+V		Nominal Ordinal	и	<u>R</u> eset

Figure 155. Configuration du niveau de mesure pour plusieurs champs

Astuce : Pour modifier les propriétés de plusieurs champs contenant des valeurs similaires (telles que 0/1), cliquez sur l'en-tête de colonne *Valeurs* afin de trier les champs en fonction de cette valeur. Maintenez la touche Maj enfoncée tout en utilisant la souris ou les touches fléchées pour sélectionner tous les champs à modifier. Cliquez ensuite sur la sélection avec le bouton droit de la souris pour modifier le niveau de mesure ou les autres attributs des champs sélectionnés.

3. Définissez le niveau de mesure pour le champ *attrition* sur **Indicateur**, puis définissez le rôle sur **Cible**. Le rôle de tous les autres champs doit être défini sur **Entrée**.

	Preview				0-[
Types Format Types Format	Annotations	alues Clear	Values	Clear All Va	ilues
Field -	Measurement	Values	Missing	Check	Role
epili S	nay	170	_	NONE	 input
loglong	Continuous	[-0.10536		None	🔪 Input
Iogtoll	🖉 Continuous	[1.74919		None	🔪 Input
🔉 logequi	🔗 Continuous	[2.73436		None	🔪 Input
logcard	🔗 Continuous	[1.01160		None	🔪 Input
logwire	Continuous	[2.70136		None	🔪 Input
Ininc	🖉 Continuous	[2.19722		None	🔪 Input
custcat	💑 Nominal	1,2,3,4		None	🔪 Input
churn	🎖 Flag	1/0		None	🔘 Target
View currer	t fields 🔘 View unu	ised field setting	IS		

Figure 156. Configuration du niveau de mesure et du rôle pour le champ attrition

4. Ajoutez au noeud type un noeud de modélisation Sélection de fonction.

L'utilisation d'un noeud Sélection de fonction vous permet de supprimer les prédicteurs ou les données qui n'apportent aucune information utile en matière de relation prédicteur/cible.

- 5. Exécutez le flux.
- 6. Ouvrez le nugget de modèle obtenu, et à partir du menu **Générer**, sélectionnez **Filtrer** pour créer un noeud Filtrer.

🖬 churn 🛛 🔀									
0	🐞 File	O Generate	Pr	eview 🐻			0		
		🖔 Generate	e Modeling	Node				<u></u>	
		Model to	Palette						
Model	Summary	Filter	N						
	- (e)		13						
		Rank	. 4	ો 1લેં					
	Rank 🚣	Field	Me	asurement		Importance	Value		
	1	🛞 tenure	🖉 Contin	nuous	*	Important	1.0	4	
-	2	Ioglong	Ontin	nuous	*	Important	1.0		
-	3	🛞 equip	💑 Nomir	al	*	Important	1.0		
-	4	🛞 longten	🔗 Contin	nuous	*	Important	1.0		
-	5	🛞 employ	🔗 Contin	nuous	*	Important	1.0		
-	6	🛞 longmon	Ontir	nuous	*	Important	1.0		
-	7	🛞 internet	💑 Nomir	al	*	Important	1.0		
-	8	🛞 equipmon	🔗 Contin	nuous	*	Important	1.0		
-	9	🛞 age	Contin	nuous	*	Important	1.0		
-	10	🛞 ebill	💑 Nomir	al	*	Important	1.0		
-	11	🛞 address	🖉 Contin	nuous	*	Important	1.0		
-	12	🛞 callcard	💑 Nomir	al	*	Important	1.0		
-	13	🛞 cardten	🖉 Contin	nuous	*	Important	1.0		
-	14	🛞 ed	- Ordina	al	*	Important	1.0		
-	15	🛞 toliten	🔗 Contin	nuous	*	Important	1.0		
-	16	🛞 custcat	💑 Nomir	al	*	Important	1.0		
-	17	🛞 voice	💑 Nomir	al	*	Important	1.0		
	18	🛞 cardmon	🖉 Contir	nuous	*	Important	1.0		
-	19	🛞 logtoll	🖉 Contin	nuous	*	Important	1.0		
	20	🛞 wireless	💑 Nomir	al	*	Important	1.0	-	
Selected fields: 27 Total fields available: 41 ★> 0.95 + <= 0.95 < 0.9									
			3 Scre	ened Fields					
	Field 😨	Measure	ment			Reason			
E	🛞 retire	💑 Nominal		Single catego	ry too	large			
	🛞 logwire	🖉 Continuou	IS	Too many mis	sing v	alues			
	🛞 logequi	Continuou	IS	Coefficient of	variat	ion below thr	eshold		
ОК	Cancel						Apply	Reset	

Figure 157. Génération d'un noeud Filtrer à partir d'un noeud Sélection de fonction

Toutes les données du fichier *telco.sav* ne sont pas utiles à la prévision de l'attrition. Vous pouvez appliquer le filtre pour ne sélectionner que les données considérées comme importantes en tant que prédicteur.

- 7. Dans la boîte de dialogue Générer un filtre, sélectionnez **Tous les champs marqués : Important** et cliquez sur **OK**.
- 8. Reliez le noeud Filtrer généré au noeud type.

Mode:	Include	C Exclude
O Selecte	d fields	
All field	s marked:	
	🗴 🔀 Important	
E	🚺 🛨 Marginal	
E	📃 💽 Unimporta	nt
🔘 Top nu	nber of fields	10 🗘
🔘 Importa	nce greater than:	0.667 🖨

Figure 158. Sélection des champs importants

- 9. Liez un noeud Audit données au noeud Filtrer généré.
 - Ouvrez le noeud Audit données, puis cliquez sur Exécuter.
- 10. Dans l'onglet Qualité du navigateur Audit données, cliquez sur la colonne % *terminé(s)* pour la trier dans l'ordre numérique croissant. Vous pouvez ainsi identifier les champs où de grandes quantités de données manquent. Dans notre exemple, le seul champ à modifier est *logtoll*, qui est complet à moins de 50 %.
- 11. Dans la colonne Attribuer une entrée manquante du champ logtoll, cliquez sur Spécifier.

Data Audit	of [28 fields] #.	2		_				
違 <u>F</u> ile 🛛 📄 E	dit 👋 <u>G</u> enerate							
Audit Quality	Annotations							
Complete fields	(%): 96.43% Co	mplete recorc	ls (%): 47.5	%		1	16 1	
Field	Measurement	Outliers	Extremes	Action	Impute Missing	Method	% Complete 🚣	Valid
📿 logtoll	Continuous	2	0	None	Never 🚿	Fixed	47.5	
决 tenure	🔗 Continuous	0	0	None	Never	Fixed	100	
💭 age	🔗 Continuous	0	0	None	Blank Values	Fixed	100	
🤔 address	🔗 Continuous	12	0	None	Null Values	Fixed	100	
🔉 income	🔗 Continuous	9	6	None	Blank & Null Value	Fixed	100	
决 ed	📶 Ordinal				Condition	Fixed	100	
🔎 employ	🔗 Continuous	8	0	None	Specify	Fixed	100	
决 equip	🎖 Flag				Never A	Fixed	100	
🔆 callcard	🎖 Flag				Never	Fixed	100	
🚯 wireless	🎖 Flag				Never	Fixed	100	
Iongmon	Continuous	18	4	None	Never	Fixed	100	
tollmon	Continuous	9	1	None	Never	Fixed	100	
🛞 equipmon	Continuous	2	0	None	Never	Fixed	100	
n 😥 cardmon	Continuous	11	3	None	Never	Fixed	100	
😥 wiremon	Continuous	8	1	None	Never	Fixed	100	
😥 longten	Continuous	20	4	None	Never	Fixed	100	
tollten	Continuous	18	2	None	Never	Fixed	100	
칮 cardten	Continuous	11	6	None	Never	Fixed	100	
🗘 voice	🖁 Flag				Never	Fixed	100	

Figure 159. Attribution des valeurs manquantes au champ logtoll

12. Dans le champ Attribuer quand, sélectionnez Valeurs nulles et non renseignées. Dans le champ Fixe en tant que, sélectionnez Moyenne et cliquez sur OK.

La sélection de **Moyenne** garantit que les valeurs attribuées n'ont pas d'impact négatif sur la moyenne de toutes les valeurs dans les données globales.

Field:	logtoll	Storage: 🛞 Real	
Impute wher	ı:	Blank & Null Values 💌	
Condition:			
Impute Metho	od:	Fixed	-
Impute Fixe	ed Values		
Fixed as:	Mean 📉		
Value:	Mean Mid-Range		

Figure 160. Sélection des paramètres d'imputation

13. Dans l'onglet Qualité du navigateur Audit données, générez le super noeud Valeurs manquantes. Pour ce faire, à partir des menus, choisissez :

Générer >	Super	noeud	des	valeurs	manq	uantes
-----------	-------	-------	-----	---------	------	--------

Audit Guarty Ar Messing Values SuperNode Outlier & Extreme SuperNode Missing Values Filter Node 5% Field Missing Values Filter Node 5% Field Missing Values Select Node 5% Iterure Binning Node 0 None Blank & Null Val Fixed 47 gage Dinning Node 0 None Never Fixed 11 address Derive Node 0 None Never Fixed 11 edup Oraph Output Never Fixed 11 edup Oraph Output Never Fixed 11 edup Nominal Never Fixed 11 edup Nominal Never Fixed 11 edup Nominal Never Fixed 11 edup Orntinuous 18 4 None Never Fixed 11 forgon Continuous 1 3 None Never	Eile 🔒 E	dit 👋 Generate							0
Outlier & Extreme SuperNode S% Field Missing Values Eliter Node S% Extremes Action Impute Missing Method % Complete IogoII Reclassity Node D None Blank & Null Val Fixed 47 IogoII Reclassity Node D None Never Fixed 41 age Dining Node Derive Node O None Never Fixed 11 address Derive Node O None Never Fixed 11 address Derive Node O None Never Fixed 11 edus Graph Qutput - - Never Fixed 11 edus Nominal - - Never Fixed 11 outliers Nominal - - Never Fixed 11 outliers 9 1 None Never Fixed 11 outliers 8 1 None Never Fixed 11 outliers 2 0 None Never Fixed 11	udit Quality	An Missing Values S	uperNode						
Sing Values Eiter Node 5% Field Missing Values Select Node Extremes Action Impute Missing Method % Complete logtoll Reclassify Node 0 None Blank & Null Val Fixed 47 logtoll Reclassify Node 0 None Never Fixed 41 age Derive Node 0 None Never Fixed 11 oddress Derive Node 0 None Never Fixed 11 income Graph Output 6 None Never Fixed 11 equip Nominal Never Fixed 11 calcrad Nominal Never Fixed 11 vireless Nominal Never Fixed 11 tollmon Continuous 18 4 None Never Fixed 11 tollmon Continuous 11 3 None Never Fixed 11 tol		Outlier & Extreme	SuperNode						
Higsing Values Select Node Extremes Action Impute Missing Method % Complete logtoll Reclassify Node 0 None Blank & Null Val Fixed 47 logtoll Reclassify Node 0 None Never Fixed 11 address Derive Node 0 None Never Fixed 11 address Derive Node 0 None Never Fixed 11 address Oraph Output - - Never Fixed 11 equip Oraph Output - - None Never Fixed 11 equip Oraph Node 0 None Never Fixed 11 equip Oraph Node 0 None Never Fixed 11 oraphon Continuous 18 4 None Never Fixed 11 termin Continuous 11 3 None Never Fixed 11 tollmon Continuous 11 3 None Never	omplete fields (^{(%):} Missing Values <u>F</u> i	lter Node	.5%					
logtoll Reclassify Node 0 None Blank & Null Val Fixed 47 logtoll Binning Node 0 None Never Fixed 11 address Derive Node 0 None Never Fixed 11 address Derive Node 0 None Never Fixed 11 address Graph Output 6 None Never Fixed 11 equip O Nominal - Never Fixed 11 equip Nominal - Never Fixed 11 vireless Nominal - Never Fixed 11 organon Continuous 18 4 None Never Fixed 11 organon Continuous 9 1 None Never Fixed 11 organon Continuous 11 3 None Never Fixed 11 organon Continuous 11 3 None Never Fixed 11 organon Continuous 11 3 None N	Field	Missing Values S	elect Node	Extremes	Action	Impute Missing	Method	% Complete 🚣	Valid
Iterure Binning Node O None Never Fixed 11 address Derive Node O None Never Fixed 11 address Derive Node O None Never Fixed 11 income O Rever Fixed 11 ed Graph Output Never Fixed 11 equip Nominal Never Fixed 11 equip Nominal Never Fixed 11 vireless Nominal Never Fixed 11 iongmon Continuous 18 4 None Never Fixed 11 iongmon Continuous 18 4 None Never Fixed 11 iongmon Continuous 11 3 None Never Fixed 11 iongmon Continuous 11 None	logtoll	Reclassify Node		0 None		Blank & Null Val	Fixed	47.5	
age Dinning Node Dinning Node address Derive Node 0 None Never Fixed 11 address Derive Node 0 None Never Fixed 11 income Graph Output 6 None Never Fixed 11 equip Graph Output 0 None Never Fixed 11 equip O Nominal Never Fixed 11 calcard Nominal Never Fixed 11 longmon Continuous 18 4 None Never Fixed 11 longmon Continuous 18 4 None Never Fixed 11 equipmon Continuous 18 4 None Never Fixed 11 iorgmon Continuous 18 4 None Never Fixed 11 iorgmon Continuous 11 3 None Never Fixed 11 iorgmon Continuous 18 1 None Never Fixed 11 iorgmon Continuous 18 None Never Fixed 11 iorgmon Continuous 18 None Never<	tenure	6		0 None		Never	Fixed	100	
address Derive Node 0 None Never Fixed 11 income Graph Output 6 None Never Fixed 11 ed Graph Output - Never Fixed 11 employ O Noninal - Never Fixed 11 equip Nominal - - Never Fixed 11 callcard Nominal - - Never Fixed 11 urless Nominal - - Never Fixed 11 longmon Continuous 18 4 None Never Fixed 11 tollmon Continuous 18 4 None Never Fixed 11 urlend Continuous 11 3 None Never Fixed 11 urlend Continuous 11 3 None Never Fixed 11 urlend Continuous 18 1 None Never Fixed 11 urlend Continuous 18 2 None Never Fixed 11 voice Nominal - Never Fixed 11 voice Nominal	age	Binning Node		0 None		Never	Fixed	100	
income Graph Output 6 None Never Fixed 11 ed Graph Node 0 None Never Fixed 11 employ Graph Node 0 None Never Fixed 11 callcard Nominal Never Fixed 11 callcard Nominal Never Fixed 11 callcard Nominal Never Fixed 11 longmon Continuous 18 4 None Never Fixed 11 longmon Continuous 9 1 None Never Fixed 11 longmon Continuous 1 3 None Never Fixed 11 cardmon Continuous 11 3 None Never Fixed 11 longten Continuous 18 1 None Never Fixed 11 longten Continuous 18 2 None Never Fixed 11 longten Continuous 11 6 None	address	Derive Node		0 None		Never	Fixed	100	
ed	income	A country of dealers of		6 None		Never	Fixed	100	
employ Graph Node 0 None Never Fixed 11 equip Nominal Never Fixed 11 callcard Nominal Never Fixed 11 callcard Nominal Never Fixed 11 longmon Continuous 18 4 None Never Fixed 11 longmon Continuous 9 1 None Never Fixed 11 equipmon Continuous 9 1 None Never Fixed 11 equipmon Continuous 11 3 None Never Fixed 11 equipmon Continuous 11 3 None Never Fixed 11 viremon Continuous 11 3 None Never Fixed 11 longten Continuous 18 2 None Never Fixed 11 longten Continuous 18 2 None Never Fixed 11 vice Nominal Never Fixed 11 vice Nominal Never Fixed 11	ed	Graph Output				Never	Fixed	100	
equip Nominal Never Fixed 11 callcard Nominal Never Fixed 11 wireless Nominal Never Fixed 11 longmon Continuous 18 4 None Never Fixed 11 longmon Continuous 9 1 None Never Fixed 11 equipmon Continuous 11 3 None Never Fixed 11 oracriton Continuous 11 3 None Never Fixed 11 longten Continuous 11 3 None Never Fixed 11 longten Continuous 18 1 None Never Fixed 11 longten Continuous 18 2 None Never Fixed 11 longten Continuous 18 2 None Never Fixed 11 longten Continuous 11 <td< td=""><td>employ</td><td>Graph Node</td><td></td><td>0 None</td><td></td><td>Never</td><td>Fixed</td><td>100</td><td></td></td<>	employ	Graph Node		0 None		Never	Fixed	100	
callcard Nominal Never Fixed 11 longmon Continuous 18 4 None Never Fixed 11 longmon Continuous 18 4 None Never Fixed 11 equipmon Continuous 9 1 None Never Fixed 11 equipmon Continuous 11 3 None Never Fixed 11 ordmon Continuous 11 3 None Never Fixed 11 longten Continuous 8 1 None Never Fixed 11 longten Continuous 18 2 None Never Fixed 11 longten Continuous 18 2 None Never Fixed 11 longten Continuous 18 2 None Never Fixed 11 cardten Continuous 11 6 None Never Fixed 11 pager Nominal Never Fixed 11 cardten	equip	Nominal				Never	Fixed	100	
wireless Nominal Never Fixed 11 longmon Continuous 18 4 None Never Fixed 11 tollmon Continuous 9 1 None Never Fixed 11 cardmon Continuous 2 0 None Never Fixed 11 cardmon Continuous 11 3 None Never Fixed 11 longten Continuous 8 1 None Never Fixed 11 longten Continuous 8 1 None Never Fixed 11 longten Continuous 18 2 None Never Fixed 11 cardten Continuous 18 2 None Never Fixed 11 cardten Continuous 11 6 None Never Fixed 11 cardten Continuous 11 6 None Never Fixed 11 cardten Nominal Never Fixed 11 cardten	callcard	💑 Nominal	822	22 229		Never	Fixed	100	
Iongmon Continuous 18 4 None Never Fixed 11 tollmon Continuous 9 1 None Never Fixed 11 equipmon Continuous 2 0 None Never Fixed 11 cardinon Continuous 11 3 None Never Fixed 11 wiremon Continuous 11 3 None Never Fixed 11 longten Continuous 20 4 None Never Fixed 11 longten Continuous 20 4 None Never Fixed 11 longten Continuous 18 2 None Never Fixed 11 ottet Continuous 18 2 None Never Fixed 11 voice Nominal Never Fixed 11 voice Nominal Never Fixed 11 retret Nominal Never Fixed 11 retret Nominal Never Fixed 11 cardten Nominal Never Fixed <t< td=""><td>wireless</td><td>💑 Nominal</td><td>842</td><td>22 -23</td><td></td><td>Never</td><td>Fixed</td><td>100</td><td></td></t<>	wireless	💑 Nominal	842	22 -23		Never	Fixed	100	
tollmon Continuous 9 1 None Never Fixed 11 equipmon Continuous 2 0 None Never Fixed 11 cardmon Continuous 11 3 None Never Fixed 11 wiremon Continuous 11 3 None Never Fixed 11 longten Continuous 20 4 None Never Fixed 11 toltten Continuous 18 2 None Never Fixed 11 cardmon Continuous 11 6 None Never Fixed 11 cardmon Continuous 11 6 None Never Fixed 11 voice Nominal Never Fixed 11 pager Nominal Never Fixed 11 callwait Nominal Never Fixed 11 callwait Nominal Never Fixed 11 callwait No	longmon	Continuous	18	4 None		Never	Fixed	100	
equipmon	tollmon	🔗 Continuous	9	1 None		Never	Fixed	100	
cardmon Continuous 11 3 None Never Fixed 11 wiremon Continuous 8 1 None Never Fixed 11 longten Continuous 20 4 None Never Fixed 11 longten Continuous 20 4 None Never Fixed 11 cardten Continuous 11 6 None Never Fixed 11 cardten Continuous 11 6 None Never Fixed 11 pager Nominal Never Fixed 11 callwait Nominal Never Fixed 11 contro Nominal Never Fixed 11 contro Nominal	equipmon	🔗 Continuous	2	0 None		Never	Fixed	100	
wiremon Continuous 8 1 None Never Fixed 11 longten Continuous 20 4 None Never Fixed 11 tolten Continuous 18 2 None Never Fixed 11 cardten Continuous 11 6 None Never Fixed 11 voice Nominal Never Fixed 11 pager Nominal Never Fixed 11 pager Nominal Never Fixed 11 callwait Never Fixed 11 callwait Never Fixed 11 callwait Never Fixed 11 callwait Never Fixed 11 contra Nominal Never Fixed 11 <td< td=""><td>cardmon</td><td>🔗 Continuous</td><td>11</td><td>3 None</td><td></td><td>Never</td><td>Fixed</td><td>100</td><td></td></td<>	cardmon	🔗 Continuous	11	3 None		Never	Fixed	100	
Iongten Continuous 20 4 None Never Fixed 11 tollten Continuous 18 2 None Never Fixed 11 cardten Continuous 11 6 None Never Fixed 11 voice Nominal Never Fixed 11 pager Nominal Never Fixed 11 internet Nominal Never Fixed 11 callwait Nominal Never Fixed 11 confer Nominal Never Fixed 11 collwait Nominal Never Fixed 11 collwait Nominal Never Fixed 11 collwait Nominal Never Fixed 11 colloga Nominal	wiremon	Continuous	8	1 None		Never	Fixed	100	
tollten Continuous 18 2 None Never Fixed 11 cardten Continuous 11 6 None Never Fixed 11 voice Nominal Never Fixed 11 pager Nominal Never Fixed 11 internet Nominal Never Fixed 11 callwait Nominal Never Fixed 11 confer Nominal Never Fixed 11 confer Nominal Never Fixed 11 confer Nominal Never Fixed 11 ebill Nominal Never Fixed 11 coldona Continuous 4 None Never Fixed 11	longten	n Continuous	20	4 None		Never	Fixed	100	
Cardten Continuous 11 6 None Never Fixed 11 voice Nominal Never Fixed 11 pager Nominal Never Fixed 11 internet Nominal Never Fixed 11 callwait Nominal Never Fixed 11 callwait Nominal Never Fixed 11 conter Nominal Never Fixed 11 conter Nominal Never Fixed 11 conter Nominal Never Fixed 11 ebill Nominal Never Fixed 11	tollten	Continuous	18	2 None		Never	Fixed	100	
voice Mominal Never Fixed 11 pager Nominal Never Fixed 11 internet Sominal Never Fixed 11 callwait Nominal Never Fixed 11 callwait Nominal Never Fixed 11 confer Nominal Never Fixed 11 ebill Nominal Never Fixed 11 colona Never Fixed 11	cardten	Continuous	11	6 None		Never	Fixed	100	
pager Nominal Never Fixed 11 Internet Nominal Never Fixed 11 callwait Nominal Never Fixed 11 callwait Nominal Never Fixed 11 confer Nominal Never Fixed 11 ebill Nominal Never Fixed 11 oldong Continuous 4 None Never Fixed 11	voice	nominal				Never	Fixed	100	
internet Nominal Never Fixed 11 callwait Nominal Never Fixed 11 confer Nominal Never Fixed 11 ebill Nominal Never Fixed 11 loalong Continuous 4 0 None Never Fixed 11	pager	nominal	8			Never	Fixed	100	
callwait Movinal Never Fixed 11 confer & Nominal Never Fixed 11 ebill & Nominal Never Fixed 11 loalona Never Fixed 11	internet	nominal	8- -			Never	Fixed	100	
confer ▲ Nominal Never Fixed 11 ebill ▲ Nominal Never Fixed 11 loalona ▲ Continuous 4 0 None Never Fixed 11	callwait	ot Nominal				Never	Fixed	100	
ebill 💑 Nominal Never Fixed 11 Ioglong 🖉 Continuous 4 0 None Never Fixed 11	confer	ot Nominal	844			Never	Fixed	100	
logiong 🖉 Continuous 4 0 None Never Fixed 11	ebill	ot Nominal	86 -2			Never	Fixed	100	
	loglong	Continuous	4	0 None		Never	Fixed	100	
Ininc S Onitinuous 9 0 None Never Fixed 11	Ininc	S Continuous	9	0 None		Never	Fixed	100	

Figure 161. Génération d'un super noeud des valeurs manquantes

Dans la boîte de dialogue Super noeud des valeurs manquantes, augmentez le paramètre **Taille d'échantillon (%)** à 50 %, puis cliquez sur **OK**.

Le super noeud apparaît dans le canevas de flux, avec l'intitulé : Attribution de valeur manquante.

14. Reliez le super noeud au noeud Filtrer.

Generate SuperNo	de for:
All fields	O Selected fields only
Sample Size (%):	50.00 ≑

Figure 162. Définition de la taille d'échantillon

- 15. Ajoutez un noeud Logistique au super noeud.
- **16**. Dans le noeud Logistique, cliquez sur l'onglet Modèle et sélectionnez la procédure **Binomial**. Dans la zone *Procédure binomiale*, sélectionnez la méthode **Ascendante**.

🔛 churn			×
		0	
Fields Model Expert	Analyze Annotations		_
Model name: 🧕 Auto 🔘	Custom		
👿 Use partitioned data			
👿 Build model for each sp	lit		
Procedure: O Multinom	ial	Binomial	
Binomial Procedure	_		
Method: Forwards			
Categorical Inputs:			
Field Name	Contrast	Base Category	
			×
👿 Include constant in equ	ation		
OK 🕨 Run Ca	ancel	Apply	Reset

Figure 163. Choix des options de modèle

- 17. Dans l'onglet Expert, sélectionnez le mode **Expert**, puis cliquez sur **Sortie**. La boîte de dialogue Sorties avancées apparaît.
- **18**. Dans la boîte de dialogue Sorties avancées, sélectionnez **A chaque étape** en tant que type *Afficher*. Sélectionnez **Historique d'itération** et **Estimations des paramètres**, puis cliquez sur **OK**.

💟 Logistic Regression: Advanced Output						
Display:) At each step	◯ At last step				
🛃 Iteration his	tory	📝 Parameter estimates				
Classificatio	on plots	Hosmer-Lemeshow goodness-of-fit				
CI for exp(E	3) (%)	95				
📃 Residual Dia	gnosis					
Outli	ers outside (std. dev.):	2.0 🖨				
🔘 All c	ases					
Classification cu	toff:	0.5 🚔				
	ОК Са	ancel Help				

Figure 164. Choix des options de sortie

Navigation dans le modèle

1. Dans le noeud Logistique, cliquez sur Exécuter pour créer le modèle.

Le nugget de modèle est ajouté à l'espace de travail du flux et également à la palette Modèles en haut à droite. Pour afficher ses détails, cliquez avec le bouton droit de la souris sur le nugget de modèle et sélectionnez **Editer** ou **Parcourir**.

L'onglet Récapitulatif affiche (entre autres) la cible et les entrées (champs prédicteurs) utilisées par le modèle. Ces champs sont ceux qui ont été réellement choisis sur la base de la méthode Ascendante, et non la liste complète soumise.

🖸 churn 🛛 🕅
🖕 🕼 File 🖏 Generate 📑 Preview 📾 🕢 🗖 🗖
Summary Advanced Settings Annotations
Collapse All 🌾 Expand All
Analysis Fields Fields Target Genute
OK Cancel Apply Reset

Figure 165. Récapitulatif du modèle avec champs cible et champs d'entrée

Les éléments affichés dans l'onglet Options avancées dépendent des options sélectionnées dans la boîte de dialogue Sorties avancées, dans le noeud Logistique. L'élément Récapitulatif du traitement des observations est systématiquement affiché. Il indique le nombre et le pourcentage d'enregistrements inclus dans l'analyse. Par ailleurs, il indique le nombre d'observations manquantes (s'il y a lieu) où un ou plusieurs des champs d'entrée ne sont pas disponibles, ainsi que toutes les observations qui n'ont pas été sélectionnées.

	churn				×				
L	File 🖏 Generate 🕞 Preview 🐻 🕢 🗖 🗖								
	Summary Advanced Settings Annotations								
	Logistic Regression								
	Unweighted Cases(a)	Process	ang summary	N	Percent				
		Include	d in Analysis	1000	100.0				
	Selected Cases	Miss	ing Cases	0	.0				
			Total	1000	100.0				
	Unselecte	d Cases		0	.0				
	Tot	al		1000	100.0				
	a. If weight is in effect, see classification table for the total number of cases.								
	Deper	ndent Var	iable Encoding	ri -					
	Original Value Internal Value								
		No	0		*				
	1				1				
0	OK				Apply Reset				

Figure 166. Récapitulatif du traitement des observations

2. Faites défiler la fenêtre vers le bas à partir de l'élément Récapitulatif du traitement des observations pour afficher la table de classification, sous Bloc 0 : le bloc de début.

La méthode Pas à pas ascendante commence par un modèle nul, c'est-à-dire un modèle sans prédicteur, qui peut servir de base à la comparaison avec le modèle final créé. Le modèle nul, par convention, donne systématiquement la valeur de prévision 0. Par conséquent, le modèle nul est précis à 72,6 %, tout simplement parce que les 726 clients n'ayant pas changé de fournisseur font l'objet d'une prévision correcte. A l'inverse, la prévision concernant les clients ayant changé de fournisseur n'est pas du tout correcte.

Cl	nurn	File 🐑 Ge	enerate ings A	() Annota	Prev tions	iew) 🛃 🕢 🗖					
b. Ir	iitial -2 Lo	g Likelihood: 11	74.394				F				
c.E cha	stimation nged by li	terminated at ite ess than .000.	eration n	umber	4 bec	ause parameter estimates					
		Cia	SSILICAT		able(a	Predicted					
				ch	urn						
		Observed		No	Yes	Percentage Correct					
		churp	No	726	0	100.0					
	Step 0	chum	Yes	274	0	.0					
		Overall Percenta	l Ige			72.6					
a. Constant is included in the model.											
	b. The c	ut value is .500									
4	Variables in the Equation										
OK	Саг	icel		OK Cancel Apply Reset							

Figure 167. Début de la table de classification supervisée - Bloc 0

3. Faites maintenant défiler la fenêtre vers le bas pour afficher la table de classification supervisée, sous Bloc 1 : Méthode = Pas à pas ascendante.

Cette table de classification supervisée affiche les résultats du modèle lors de l'ajout d'un prédicteur à chacune des étapes. Dès la première étape (alors qu'un seul prédicteur a été ajouté), le modèle a permis d'augmenter l'exactitude de la prévision d'attrition de 0,0 % à 29,9 %

CI	Churn 🛛								
Sum	Summary Advanced Settings Annotations								
1	Predicted								
				ch	urn	Percentage Correct			
		Observed		No	Yes				
		churn	No	668	58	92.0			
	Step 1		Yes	192	82	29.9			
		Overal Percenta	l Ige			75.0			
		churp	No	657	69	90.5			
	Step 2	chun	Yes	160	114	41.6			
		Overall Percentage				77.1			
		churp	No	661	65	91.0			
	Step 3		Yes	153	121	44.2			
1									
OK Cancel Apply Reset									

Figure 168. Table de classification supervisée - Bloc 1

4. Faites défiler la table de classification supervisée jusqu'en bas.

La table de classification supervisée montre que la dernière étape est l'étape 8. A ce stade, l'algorithme a décidé qu'il est inutile d'ajouter d'autres prédicteurs au modèle. Bien que l'exactitude des clients n'ayant pas changé de fournisseur ait diminué quelque peu jusqu'à 91,2%, l'exactitude de la prévision des clients ayant changé de fournisseur a augmenté de la valeur d'origine 0 %, à 47,1 %. Cela représente une amélioration significative par rapport au modèle nul d'origine sans prédicteur.

mmary Ad	File 🖏 Gel vanced Setti	nerate ngs A	() Annota	Pre tions	view			0)[=	
otop v	Overall Percentag	je						78.7	
		No	657	69				90.5	
Step 7	churn	Yes	144	130				47.4	
	Overall Percentage					78.7			
	ahura	No	662	64		91		91.2	Í
Step 8	cnurn	Yes	145	129				47.1	
	Overall Percentag	je						79.1	
a. The cu	a. The cut value is .500								
	Varia	ables i	n the	Equat	ion	_	_		a l
B S.E. Wald dt				df	Sig.	Exp(B)			
Step 1(a)	tenure	046	.004	123	3.346	1	.000	.955	
Step 1(a) Constant 462 136 11 574 1 001 1 587									
K Cancel Apply Reset									

Figure 169. Table de classification supervisée - Bloc 1

Pour un client qui souhaite réduire l'attrition, le fait de pouvoir la réduire presque de moitié est une étape majeure de protection des flux de revenus.

Remarque : Cet exemple montre également que le fait de prendre le pourcentage global comme guide d'exactitude d'un modèle peut, dans certains cas, être trompeur. Le modèle nul d'origine avait une exactitude globale de 72,6 %, alors que le modèle de prévision final présente une exactitude globale de 79.1%. Toutefois, comme nous l'avons vu, l'exactitude des prévisions réelles par catégorie était très différente.

Pour évaluer le niveau d'adéquation du modèle aux données, divers diagnostics sont disponibles dans la boîte de dialogue Sorties avancées lorsque vous créez le modèle. Des explications sur les fondements mathématiques des méthodes de modélisation utilisées dans IBM SPSS Modeler sont présentées dans le *guide des algorithmes d'IBM SPSS Modeler*, disponible dans le répertoire *Documentation* du disque d'installation.

Sachez également que ces résultats sont établis uniquement d'après les données d'apprentissage. Pour évaluer comment le modèle peut se généraliser à d'autres données dans le monde réel, vous devez utiliser un noeud Partitionner destiné à contenir un sous-ensemble d'enregistrements à des fins de test et de validation.

Chapitre 14. Prévision de l'utilisation de la bande passante (Séries temporelles)

Prévision avec le noeud Séries temporelles

Un analyste pour un fournisseur large bande national doit établir des prévisions sur les abonnements des utilisateurs afin de prédire l'utilisation de la bande passante. Les prévisions sont requises pour chacun des marchés locaux qui constituent la base nationale de l'abonné. Vous utiliserez la modélisation des séries temporelles afin de produire des prévisions sur un sous-ensemble des marchés locaux pour les trois prochains mois. Un second exemple montre comment convertir des données source si leur format ne permet pas de les intégrer au noeud Séries temporelles.

Ces exemples utilisent le flux intitulé *broadband_create_models.str*, qui fait référence au fichier de données *broadband_1.sav*. Ces fichiers sont disponibles dans le dossier *Demos* de toutes les installations d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *broadband_create_models.str* se trouve dans le dossier des *flux*.

Le dernier exemple montre comment appliquer les modèles enregistrés à un jeu de données mis à jour afin de prolonger les prévisions d'une nouvelle période de trois mois.

Dans IBM SPSS Modeler, vous pouvez générer plusieurs modèles de séries temporelles simultanément. Le fichier source que vous utiliserez comporte les séries temporelles de 85 marchés différents. Toutefois, pour des raisons de simplicité, vous ne modélisez que cinq de ces marchés, ainsi que le total de tous les marchés.

Le fichier de données *broadband_1.sav* comporte les données d'utilisation mensuelle de chacun des 85 marchés locaux. Dans le cadre de cet exemple, seules les cinq premières séries seront utilisées ; un modèle distinct sera créé pour chacune de ces cinq séries, ainsi que pour un total.

En outre, le fichier comprend un champ de date qui indique le mois et l'année de chaque enregistrement. Ce champ sera utilisé pour l'étiquetage des enregistrements. Dans IBM SPSS Modeler, le champ de date est lu en tant que chaîne. Cependant, pour utiliser ce champ dans IBM SPSS Modeler, vous convertirez le type de stockage au format de date numérique à l'aide d'un noeud Remplacer.



Figure 170. Flux d'échantillons illustrant la modélisation des séries temporelles

Le noeud Séries temporelles requiert que chaque série se trouve dans une colonne distincte, avec une ligne par intervalle. IBM SPSS Modeler met à votre disposition des méthodes qui permettent, au besoin, de convertir les données pour qu'elles soient compatibles avec ce format.

😂 <u>F</u> ile	📄 Edit	🏷 Gene	rate 🚺					•	>)>
Table	Annotations								
	Market_1	Market_2	Market_3	Market_4	Market_5	Market_6	Market_7	Market_8	Mar
1	3750	11489	11659	4571	2205	5488	6144	2363	5041
2	3846	11984	12228	4825	2301	5672	6390	2404	5160
3	3894	12266	12897	5041	2352	5802	6670	2469	5232
4	4010	12801	13716	5211	2490	5899	6929	2574	540:
5	4147	13291	14647	5383	2534	6017	7312	2654	554:
6	4335	13828	15419	5496	2664	6137	7493	2699	577:
7	4554	14273	16108	5747	2738	6250	7702	2786	5904
8	4744	14664	16958	5885	2754	6439	7965	2847	603:
9	4885	15130	17642	6053	2874	6701	8107	2967	6150
10	5020	15851	18453	6229	2975	6957	8366	3099	634:
11	5208	16509	19181	6320	3042	7111	8684	3195	663:
12	5379	17225	19885	6499	3095	7275	8997	3341	6765
13	5574	18173	20565	6593	3199	7380	9326	3376	7021
14	5828	19287	21155	6680	3207	7633	9543	3443	7339
15	5942	20171	21655	6757	3298	7985	9673	3617	749
16	6139	21379	21964	6804	3387	8236	9934	3732	7716
17	6244	22067	22756	6915	3450	8464	10211	3831	7946
18	6274	23074	23464	7035	3528	8575	10440	3886	829:
19	6347	23729	24324	7151	3546	8817	10763	3938	8584
20	6399	24803	25351	7304	3604	9041	11012	3953	871
	4								1

Figure 171. Données d'abonnement mensuelles pour les marchés locaux large bande

Création du flux

- 1. Créez un nouveau flux, puis ajoutez un noeud source Statistics pointant vers le fichier broadband_1.sav.
- 2. Utilisez un noeud Filtrer pour filtrer les champs *Market_6* à *Market_85*, ainsi que les champs *MONTH_* et *YEAR_*, afin de simplifier le modèle.

Astuce : Pour sélectionner simultanément plusieurs champs adjacents, cliquez sur le champ *Market_6*, maintenez le bouton gauche de la souris enfoncé et faites glisser le curseur vers le bas jusqu'au champ *Market_85*. Les champs sélectionnés sont surlignés en bleu. Pour ajouter les autres champs, maintenez la touche Ctrl enfoncée et cliquez sur les champs *MONTH_* et *YEAR_*.

💟 Filter		X
Filter Annotations		
7.	Fields:	89 in, 82 filtered, 0 renamed, 7 out
Field -	Filter	Field
Market 80	-×->	Market 80 📥
Market_81	×>	Market_81
Market_82	→×→	Market_82
Market_83	→×→	Market_83
Market_84	_ × →	Market_84
Market_85	_ × →	Market_85
Total	\rightarrow	Total
YEAR_	_ × →	YEAR_
MONTH_	_ ★ →	MONTH_
DATE_	\rightarrow	DATE_
View current fields View Cancel	ew unused fiel	d settings

Figure 172. Simplification du modèle

Analyse des données

Il s'avère toujours utile d'examiner la nature de vos données avant de construire un modèle. Les données présentent-elles des variations saisonnières ? Bien que Expert Modeler puisse rechercher automatiquement le meilleur modèle saisonnier ou non saisonnier pour chaque série, vous pouvez souvent obtenir des résultats plus rapidement en limitant la recherche aux modèles non saisonniers en l'absence d'effets saisonniers dans les données. Sans examiner les données de chacun des marchés locaux, nous pouvons obtenir une image approximative de la présence ou de l'absence d'effets saisonniers en traçant le nombre total d'abonnés pour l'ensemble des cinq marchés.

[9]
Plot Appearance Output Annotations
Plot: Selected series Selected Time Series models
Total
Series:
X axis label: Default Custom
Display series in separate panels
Display: 🔽 Line
Point
Smoother
Limit records Maximum number of records to plot: 2000
OK Run Cancel Apply Reset

Figure 173. Traçage du nombre total d'abonnés

- 1. Dans la palette Graphiques, reliez un noeud Tracé horaire au noeud Filtrer.
- 2. Ajoutez le champ *Total* à la liste Série.
- 3. Désélectionnez les cases Afficher les séries dans des panneaux distincts et Normaliser.
- 4. Cliquez sur Exécuter.



Figure 174. Tracé horaire du champ Total

La série présente une tendance ascendante très régulière, sans aucun signe de variations saisonnières. Il peut exister des séries spécifiques ayant des effets saisonniers, mais il semble que ceux-ci ne soient pas une caractéristique générale des données.

Bien entendu, vous devez examiner chacune des séries avant d'exclure les modèles saisonniers. Vous pouvez ensuite mettre à part les séries qui présentent des effets saisonniers et les modéliser séparément.

IBM SPSS Modeler facilite le traçage de plusieurs séries simultanément.



Figure 175. Traçage de plusieurs séries temporelles

- 5. Rouvrez le noeud Tracé horaire.
- 6. Supprimez le champ Total de la liste Série (sélectionnez-le, puis cliquez sur le bouton X rouge).
- 7. Ajoutez les champs *Market_1* à *Market_5* à la liste.
- 8. Cliquez sur Exécuter.



Figure 176. Tracé horaire de plusieurs champs

L'inspection de chacun des marchés met en évidence une tendance ascendante régulière dans chaque cas. Bien que certains marchés soient un peu plus imprévisibles que d'autres, rien n'atteste la présence d'effets saisonniers.

Définition des dates

Maintenant, vous devez attribuer le type de stockage de format Date au champ DATE_.

- 1. Reliez un noeud Remplacer au noeud Filtrer.
- 2. Ouvrez le noeud Remplacer, puis cliquez sur le bouton de sélection de champ.
- 3. Sélectionnez le champ DATE_ afin de l'ajouter à la zone Renseigner les champs.
- 4. Attribuez à la condition Remplacer la valeur Toujours.
- 5. Attribuez à l'option **Remplacer par** la valeur **to_date(DATE_)**.

🚱 Filler	×
Preview)	0
Settings Annotations	
Fill in fields:	
& DATE_	×
Replace: Always	
Condition:	
@BLANK(@FIELD)	
Replace with:	
to_date(DATE_)	
OK Cancel	Apply Reset

Figure 177. Définition du type de stockage de date

Modifiez le format de date par défaut afin qu'il corresponde au format du champ Date. Cette opération permet de convertir le champ Date correctement.

- 6. Dans le menu, choisissez l'option **Outils > Propriétés du flux > Options** pour afficher la boîte de dialogue des options de flux.
- 7. Sélectionnez la sous-fenêtre **Date/Heure** et attribuez à l'option **Format de date** par défaut la valeur **MOIS AAAA**.

broadband_create_models								
						0		
Options Messages	Parameters Deployment	Execution	Globals	Search	Comments	Annotations		
Select a setting:								
General	These settings control th	ne format of	date and	time exp	ressions in t	he current		
Date/Time	stream. Click Save As D streams	efault to us	e these s	ettings a	s the default	for all your		
Number formats								
Optimization	Import date/time as:	Oate/Ti	me © St	ri <u>ng</u>				
Logging and Status	Date format:	MON YYY	Y	-				
Layout	Time format:	HH:MM:S	S	-	Rollover days	s/mins		
Geospatial	Date baseline (1st Jan):	1900						
Cocopulat	2 digit dates start from:	1030						
		0						
	nme <u>z</u> one.	Server						
L]				0			
					20	ive AS Delault		
OK Cancel					A	pply <u>R</u> eset		

Figure 178. Définition du format de date

Définition des cibles

- 1. Ajoutez un noeud type afin de définir le rôle sur **Aucun** pour le champ *DATE_*. Définissez le rôle sur **Cible** pour tous les autres champs (les champs *Market_n* ainsi que le champ *Total*).
- 2. Cliquez sur le bouton Lire les valeurs pour remplir la colonne Valeurs.

Type								
4- 00	🍋 🔰 🕨 🕨	es Clear V	alues	Clear All Value	es			
Field -	Measurement	Values	Missing	Check	Role			
🚫 Market_1	🔗 Continuous	[3750,117		None	O Target			
🚫 Market_2	Continuous	[11489,53		None	O Target			
🚫 Market_3	Continuous	[11659,60		None	Target			
🔆 Market_4	Continuous	[4571,179		None	Target			
🚫 Market_5	Continuous	[2205,6611]		None	🔘 Target			
🚫 Total	Continuous	[536413,2		None	🔘 Target			
DATE_	Continuous	[1999-01		None	O None			
View current OK Cancel	fields 🔘 View unuse	d field settings	į		Apply Reset			

Figure 179. Définition du rôle pour plusieurs champs

Définition des intervalles de temps

- 1. A partir de la palette Modélisation, ajoutez un noeud Séries temporelles au flux et reliez-le au noeud Type.
- 2. Dans l'onglet Spécifications des données de la sous-fenêtre Observations, sélectionnez DATE_ comme champ Date/Heure.
- 3. Sélectionnez Mois comme Intervalle de temps.

📀 6 fields					
Fields Data Specifications	Build Options	Model Options	Annotations		
<u>S</u> elect an item:					
Observations	💿 <u>O</u> bservat	tions are specifie	ed by a date/time field		
Time Interval	Date/time field:				
Aggregation and Distribution	🔗 DAT	Έ_	_]		
Missing Value Handling	T <u>i</u> me in	terval: Months	*		

Figure 180. Définition de l'intervalle de temps

- 4. Dans l'onglet Options de modèle, cochez la case Etendre les enregistrements dans le futur.
- 5. Paramétrez la valeur sur 3.

🕐 6 fiel	ds					
_			Y			
Fields	Data	Specifications	Build Options	Model Options	Annotations	
Model n	iame:	Auto O C	usto <u>m</u>			
			or of			
<u>C</u> onfid	lence l	imit width (%):	95.0 🔽			
Co.	<u>n</u> tinue	estimation usi	ng existing mo	del(s)		
🔲 <u>B</u> ui	ild sco	ring model only	у			
Forec	east —					
VE:	<u>x</u> tend i	records into the	e future			3 韋



Création du modèle

- 1. Dans le noeud Séries temporelles, choisissez l'onglet Champs. Dans la liste **Champs**, sélectionnez les cinq marchés et copiez-les dans les listes **Cibles** et **Entrées candidates**. Par ailleurs, sélectionnez et copiez le champ Total dans la liste **Cibles**.
- Choisissez l'onglet Options de génération et, dans la sous-fenêtre Général, assurez-vous que la Méthode Expert Modeler est sélectionnée avec tous les paramètres par défaut. Ainsi, Expert Modeler peut décider du meilleur modèle à utiliser pour chaque série temporelle. Cliquez sur Exécuter.

😵 6 fields	
Fields Data S	pecifications Build Options Model Options Annotations
elect an item:	
General	Method: Expert Modeler 🛛 👻
Output	Model Type All models Exponential smoothing models only ARIMA models only Expert Modeler considers seasonal models
	Outliers

Figure 182. Choix d'Expert Modeler pour les séries temporelles

- 3. Reliez le nugget de modèle de séries temporelles au noeud Séries temporelles.
- 4. Reliez un noeud Table au nugget de modèle de séries temporelles et cliquez sur Exécuter.



Figure 183. Flux d'échantillons illustrant la modélisation des séries temporelles

Trois nouvelles lignes (61 à 63) sont désormais ajoutées aux données d'origine. Il s'agit des lignes relatives à la période de prévision, en l'occurrence janvier à mars 2004.

Plusieurs nouvelles colonnes sont également présentes maintenant ; les colonnes TS- sont ajoutées par le noeud Intervalles de temps. Les colonnes indiquent les informations suivantes pour chaque ligne (c'est-à-dire, pour chaque intervalle dans les séries temporelles) :

Colonne	Description
\$TS-nom_colonne	Données de modèle générées pour chaque colonne des données d'origine.
\$TSLCI-nom_colonne	Valeur inférieure de l'intervalle de confiance pour chaque colonne de données du modèle généré.
\$TSUCI-nom_colonne	Valeur supérieure de l'intervalle de confiance pour chaque colonne de données du modèle généré.
\$TS-Total	Total des valeurs de \$TS-nom_colonne pour cette ligne.
\$TSLCI-Total	Total des valeurs de \$TSLCI-nom_colonne pour cette ligne.
\$TSUCI-Total	Total des valeurs de \$TSUCI- <i>nom_colonne</i> pour cette ligne.

Les colonnes les plus significatives pour la prévision sont les colonnes *\$TS-Market_n*, *\$TSLCI-Market_n* et *\$TSUCI-Market_n*. En particulier, dans les lignes 61 à 63, ces colonnes contiennent les données prévisionnelles sur les abonnements des utilisateurs et les intervalles de confiance de chacun des marchés locaux.

Examen du modèle

1. Double-cliquez sur le nugget de modèle Séries temporelles et sélectionnez l'onglet Sortie afin d'afficher les données relatives aux modèles créés pour chacun des marchés.

Output Settings Summary Ann	notations				
H 🕒 🖉 🖉 🖝 🖜 🌩 🔶	+ - 🗋 🗐 🗃	🔁 🗟			
Output 🖌 🖉 🕹 🕹 Output 🖉 🕹 🎽 Output 🖉 Time Series	Target: Mai	′ket_2			
Temporal Informatio		Mo			
Title	Model Building Method				
🖓 🖓 Parameter Esti	Number of Predi	ctors			
E Market_1	Model Fit	MSE			
Title		RMSE RMSPE			
Model Informat					
Parameter Esti		MAE			
📮 🖪 Market_3		MAPE			
← 🕼 Parameter Esti ➡ Market_3 Title Model Informat		MAYAE			
Correlogram		WAXAPE			
Parameter Esti		AIC			
Title		BIC			
		R Square			
\overline 👬 Correlogram		Stationary R Square			
Parameter Esti	Ljung-Box Q(#)	Statistic			
🖨 😼 Market_5	Ljung Dox a(#)	df			
Title		Significance			
Im Model Informat		orginitatice			

Figure 184. Modèles de séries temporelles générés pour les marchés

Dans la colonne Sortie de gauche, sélectionnez les **Informations sur le modèle** pour l'un des marchés. La ligne **Nombre de prédicteurs** indique le nombre de champs utilisés comme prédicteurs pour chaque cible, soit zéro en l'occurrence.

Les lignes restantes des tables **Informations sur le modèle** montrent différentes mesures de qualité d'ajustement pour chaque modèle. La valeur **R carré stationnaire** estime dans quelle proportion le modèle explique la variation totale de la série. Plus la valeur est élevée (avec un maximum de 1,0), meilleur est l'ajustement du modèle.

Les lignes **Statistiques Q(#)**, **df** et **Signification** sont associées à la statistique Ljung-Box qui est un test de l'aspect aléatoire des erreurs résiduelles du modèle ; plus les erreurs sont aléatoires, meilleur peut être le modèle. **Q(#)** est la statistique Ljung-Box elle-même alors que **df** (degrés de liberté) indique le nombre de paramètres de modèles qui sont libres de varier lors de l'estimation d'une cible spécifique.

La ligne **Signification** affiche la valeur de signification de la statistique Ljung-Box, qui indique si le modèle est correctement spécifié. Une valeur de signification inférieure à 0,05 indique que les erreurs résiduelles ne sont pas aléatoires ce qui implique l'existence, dans la série observée, d'une structure inexpliquée par le modèle.

Si l'on prend en compte les valeurs **R carré stationnaire** et **Signification**, les modèles choisi par Expert Modeler pour *Market_3* et *Market_4* sont acceptables. Les valeurs **Signification** pour *Market_1*, *Market_2* et *Market_5* sont toutes inférieures à 0,05, ce qui signifie qu'une expérimentation avec des modèles mieux adaptés pour ces marchés pourrait s'avérer nécessaire.

Une série de mesures de qualité d'ajustement supplémentaires apparaît. La valeur **R carré** donne une estimation de la variation totale de la série temporelle qui peut être expliquée par le modèle. Avec une valeur maximum de 1.0 pour cette statistique, nos modèles sont acceptables à cet égard.

RMSE est l'erreur moyenne quadratique, une mesure de la différence entre les valeurs réelles d'une série et les valeurs prédites par le modèle et est exprimée dans la même unité que celle utilisée pour la série elle-même. En tant que mesure d'une erreur, cette valeur doit être aussi basse que possible. A première vue, il semble que les modèles de *Market_2* et *Market_3*, tout en étant acceptables selon les statistiques produites jusque là, sont moins réussis que ceux des trois autres marchés.

Ces mesures de qualité d'ajustement supplémentaires comprennent l'erreur absolue moyenne en pourcentage (MAPE) et sa valeur maximale (MAXAPE). L'erreur absolue en pourcentage mesure la proportion de la variation d'une série cible par rapport à son niveau prévu par le modèle et elle est exprimée en valeur de pourcentage. En examinant la moyenne et le maximum sur l'ensemble des modèles, vous pouvez obtenir une indication quant à l'incertitude inhérente à vos prévisions.

La valeur MAPE indique que tous les modèles affichent une incertitude moyenne d'environ 1 %, ce qui est très bas. La valeur MAXAPE affiche l'erreur maximale absolue en pourcentage et permet d'imaginer le pire des scénarios pour vos prévisions. Elle indique que l'erreur maximale en pourcentage pour la plupart des modèles se situe approximativement entre 1,8 et 3,7 % qui, de nouveau, sont des valeurs très basses ; seul *Market_4* est supérieur, avec une valeur plus près de 7 %.

La valeur **MAE** (erreur absolue moyenne) indique la moyenne des valeurs absolues des erreurs de prévision. Comme la valeur RMSE, elle est exprimée dans les mêmes unités que celles utilisées pour la série même. **MAXAE** indique l'erreur maximale de prévision dans les mêmes unités et le pire scénario de prévisions possible.

Bien que ces valeurs absolues soient intéressantes, il est préférable d'observer les valeurs des erreurs en pourcentage (MAPE et MAXAPE) car les séries cibles représentent les quantités d'abonnés appartenant à des marchés de taille variable.

Les valeurs MAPE et MAXAPE représentent-elles un niveau d'incertitude acceptable pour les modèles ? Elles sont certainement très basses. Dans ce genre de situation, les facteurs relatifs au développement de l'entreprise entrent en jeu car le niveau de risque acceptable change d'un problème à l'autre. Nous supposons que les qualités d'ajustement des statistiques se trouvent dans des limites acceptables et examinons à présent les erreurs résiduelles.

L'examen des valeurs de la fonction des autocorrélations (ACF) et de la fonction des autocorrélations partielles (PACF) pour les résidus des modèles donne un aperçu plus quantitatif sur les modèles que la simple consultation des statistiques de qualité d'ajustement.

Un modèle de série temporelle bien défini capturera toutes les variations non-aléatoires, y compris l'effet des saisons, la tendance, la nature cyclique et les autres facteurs importants. Le cas échéant, aucune erreur ne devra être corrélée avec elle-même (autocorrélation). Toute structure significative dans l'une de ces fonctions d'autocorrélation impliquerait que le modèle sous-jacent soit incomplet.

2. Pour le quatrième marché, dans la colonne de gauche, cliquez sur **Corrélogramme** pour afficher les valeurs de la fonction d'autocorrélation (ACF) et de la fonction d'autocorrélation partielle (PACF) pour les erreurs résiduelles du modèle.



Figure 185. Valeurs des fonctions d'autocorrélation (ACF) et d'autocorrélation partielle (PACF) pour le quatrième marché

Dans ces graphiques, les valeurs d'origine de la variable d'erreur ont été découpées en 24 périodes et comparées avec la valeur d'origine afin de savoir s'il existera des corrélations. Pour que le modèle soit acceptable, aucune des barres dans le graphique supérieur (ACF) ne doit dépasser la zone grisée, que ce soit dans les plus (vers le haut) ou dans les moins (vers le bas).

Si cela se produisait, il vous faudrait vérifier le graphique inférieur (PACF) pour savoir si la structure y est confirmée. Le graphique PACF examine les corrélations après avoir contrôlé les valeurs des séries aux points temporels intermédiaires.

Les valeurs de *Market_4* étant toutes à l'intérieur de la zone grisée, nous pouvons continuer et vérifier les valeurs des autres marchés.

3. Cliquez sur le **Corrélogramme** pour chacun des autres marchés et les totaux.

Les valeurs des autres marchés montrent toutes des valeurs en dehors de la zone ombrée, ce qui confirme ce que nous suspections auparavant à partir de leurs valeurs **Signification**. Il nous faudra utiliser d'autres modèles pour ces marchés afin de savoir s'il en existe un plus adapté, mais pour le reste de cet exemple, nous nous concentrerons sur ce que nous pouvons encore apprendre du modèle *Market_4*.

- 4. Dans la palette Graphiques, reliez un noeud Tracé horaire au nugget de modèle Séries temporelles.
- 5. Dans l'onglet Tracé, désélectionnez la case à cocher Afficher les séries dans des panneaux distincts.
- **6**. Dans la liste **Série**, cliquez sur le bouton de sélection de champ, sélectionnez les champs *Market_4* et *\$TS-Market_4*, puis cliquez sur **OK** pour les ajouter à la liste.
- 7. Cliquez sur **Exécuter** pour afficher un graphique linéaire des données réelles et prévisionnelles du premier marché local.

Plot Appearance Output Annotations
Plot: Selected series Selected Time Series models
Series: Arriset_4
X axis label: 💿 Default 🛇 Custom
🔲 Display series in separate panels 🛛 🗹 Normalize
Display: 📝 Line
Point
🥅 Smoother
Limit records Maximum number of records to plot: 2000
OK P Run Cancel Apply Reset

Figure 186. Sélection des champs à tracer

Comme vous pouvez le constater, la ligne prévisionnelle (*\$TS-Market_4*) s'étend au-delà de la fin des données réelles. Vous disposez désormais d'une prévision de la demande pour les trois prochains mois dans ce marché.

Dans le graphique, les lignes des données réelles et prévisionnelles sont très proches l'une de l'autre sur la totalité de la série temporelle, ce qui indique que ce modèle est fiable pour cette série.



Figure 187. Tracé horaire des données réelles et prévisionnelles pour Market_4

Enregistrez le modèle dans un fichier pour l'utiliser ultérieurement dans un autre exemple :

- 8. Cliquez sur OK pour fermer le graphique actuel.
- 9. Ouvrez le nugget de modèle Séries temporelles.
- 10. Choisissez l'option Fichier > Enregistrer le noeud et spécifiez l'emplacement du fichier.
- 11. Cliquez sur Enregistrer.

Vous disposez d'un modèle fiable pour ce marché en particulier, mais quelle est la marge d'erreur de la prévision ? Vous pouvez obtenir une indication en examinant l'intervalle de confiance.

- 12. Cliquez deux fois sur le dernier noeud Tracé horaire dans le flux (celui nommé Market_4 \$TS-Market_4) pour rouvrir sa boîte de dialogue.
- **13**. Cliquez sur le bouton de sélection de champ et ajoutez les champs *\$TSLCI-Market_4* et *\$TSUCI-Market_4* à la liste **Série**.
- 14. Cliquez sur Exécuter.



Figure 188. Ajout d'autres champs à tracer

Vous obtenez le même graphique qu'auparavant avec, en plus, les limites supérieure (*\$TSUCI*) et inférieure (*\$TSLCI*) de l'intervalle de confiance.

Comme vous pouvez le constater, les limites de l'intervalle de confiance divergent sur l'ensemble de la période de prévision, ce qui traduit une incertitude croissante dès lors que la prévision porte sur une période plus longue.

Toutefois, au terme de chaque période, vous disposez, en l'occurrence, d'un mois supplémentaire de données d'utilisation réelles sur lesquelles vous pouvez baser vos prévisions. Vous pouvez lire les nouvelles données du flux et réappliquer le modèle, puisque vous savez que celui-ci est fiable. Pour plus d'informations, voir la rubrique «Réapplication d'un modèle de séries temporelles», à la page 171.



Figure 189. Tracé horaire comportant l'intervalle de confiance

Récapitulatif

Vous avez appris à utiliser Expert Modeler afin de produire des prévisions pour plusieurs séries temporelles et avez enregistré les modèles obtenus dans un fichier externe.

Dans l'exemple suivant, vous allez apprendre à convertir les séries temporelles non standard dans un format permettant de les intégrer à un noeud Séries temporelles.

Réapplication d'un modèle de séries temporelles

Cet exemple applique les modèles de séries temporelles issus du premier exemple de séries temporelles, mais il peut être utilisé indépendamment. Pour plus d'informations, voir la rubrique «Prévision avec le noeud Séries temporelles», à la page 153.

Comme dans le scénario d'origine, un analyste pour un fournisseur de large bande national doit établir des prévisions mensuelles sur les abonnements des utilisateurs pour chaque marché d'une série de marchés locaux, afin de prédire les exigences en matière de bande passante. Vous avez déjà utilisé Expert Modeler pour créer des modèles et pour établir une prévision portant sur une période de trois mois.

Votre entrepôt de données ayant été mis à jour avec les données réelles correspondant à la période prévisionnelle d'origine, vous souhaitez utiliser ces données pour étendre l'horizon de prévision d'une nouvelle période de trois mois.

Cet exemple utilise le flux intitulé *broadband_apply_models.str*, qui référence le fichier de données *broadband_2.sav*. Ces fichiers sont disponibles dans le dossier *Demos* de toutes les installations d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *broadband_apply_models.str* se trouve dans le dossier des *flux*.

Récupération du flux

Dans cet exemple, vous allez recréer un noeud Séries temporelles à partir du modèle de séries temporelles enregistré dans le premier exemple. Ne vous inquiétez pas si vous ne disposez d'aucun modèle enregistré ; le dossier *Demos* en contient déjà un.

1. Ouvrez le flux broadband_apply_models.str à partir du dossier des flux du dossier Demos.



Figure 190. Ouverture du flux

Les données mensuelles mises à jour sont collectées dans le fichier broadband_2.sav.

2. Reliez un noeud Table au noeud source du fichier IBM SPSS Statistics, ouvrez le noeud Table et cliquez sur **Exécuter**.

Remarque : Le fichier de données a été mis à jour avec les données réelles relatives aux ventes réalisées de janvier à mars 2004, dans les lignes 61 à 63.

- Can									-
File File	📄 Edit 🛛 🖑) <u>G</u> enerate			l			0	×
Table A	nnotations								
	1 Market_82	Market_83	Market_84	Market_85	Total	YEAR_	MONTH_	DATE_	
44	58820	20482	14326	16935	17917	2002	8	AUG 2002	4
45	60119	21211	14349	17179	18249	2002	9	SEP 2002	L
46	61320	21893	14333	17601	18601	2002	10	OCT 2002	
47	63099	22471	14229	17816	18945	2002	11	NOV 2002	
48	64687	23112	14514	17937	19343	2002	12	DEC 2002	
49	65518	23686	14856	18003	19752	2003	1	JAN 2003	
50	65570	24669	15182	17875	20148	2003	2	FEB 2003	
51	66567	25469	15709	18214	20540	2003	3	MAR 2003	
52	67527	25868	16155	18557	20922	2003	4	APR 2003	
53	67724	26284	16521	19190	21300	2003	5	MAY 2003	
54	68644	26468	16567	19938	21669	2003	6	JUN 2003	
55	69878	26781	16618	20876	22004	2003	7	JUL 2003	
56	71538	27566	16553	21514	22398	2003	8	AUG 2003	
57	73162	28164	16597	21779	22773	2003	9	SEP 2003	
58	74167	28693	16669	22266	23160	2003	10	OCT 2003	
59	76036	28922	16748	22559	23616	2003	11	NOV 2003	
60	76630	29811	16798	23018	24067	2003	12	DEC 2003	
61	79002	30034	17122	23160	24509	2004	1	JAN 2004	
62	81123	30091	17581	23698	24968	2004	2	FEB 2004	
63	83909	30162	17894	24355	25383	2004	3	MAR 2004	۲
	4								

Figure 191. Données de vente mises à jour
Extraction du modèle enregistré

 Dans le menu IBM SPSS Modeler, choisissez l'option Insérer > Noeud depuis le fichier, puis sélectionnez le fichier *TSmodel.nod* dans le dossier *Demos* (ou utilisez le modèle de séries temporelles que vous avez enregistré dans le premier exemple de séries temporelles).

Ce fichier contient les modèles de séries temporelles issus de l'exemple précédent. L'opération d'insertion place le nugget de modèle Séries temporelles correspondant sur l'espace de travail.



Figure 192. Ajout du nugget de modèle

Génération d'un noeud de modélisation

1. Ouvrez le nugget de modèle Séries temporelles et sélectionnez l'option **Générer > Générer le noeud de modélisation**.

Cette opération place un noeud de modélisation Séries temporelles sur l'espace de travail.



Figure 193. Génération d'un noeud de modélisation à partir du nugget de modèle

Génération d'un nouveau modèle

 Fermez le nugget de modèle Séries temporelles puis supprimez-le de l'espace de travail. L'ancien modèle a été construit à partir de 60 lignes de données. Vous devez générer un nouveau modèle à partir des données de vente mises à jour (63 lignes). 2. Reliez au flux le noeud génération Séries temporelles nouvellement généré.



Figure 194. Association du noeud modélisation au flux

- 3. Ouvrez le noeud Séries temporelles.
- 4. Dans l'onglet **Options de modèle**, vérifiez que l'option **Poursuivre l'estimation à l'aide des modèles existants** est activée.



Figure 195. Réutilisation des paramètres stockés pour la modélisation des séries chronologiques

- 5. Vérifiez que l'option Etendre les enregistrements dans le futur a pour valeur 3.
- 6. Cliquez sur **Exécuter** pour mettre un nouveau nugget de modèle sur l'espace de travail et dans la palette Modèles.

Examen du nouveau modèle

- 1. Reliez un noeud Table au nouveau noeud de modèle Séries temporelles sur l'espace de travail.
- 2. Ouvrez le noeud Table, puis cliquez sur Exécuter.

Le nouveau modèle effectue toujours des prévisions sur trois mois car vous réutilisez les paramètres stockés. Néanmoins, cette fois, la prévision porte sur la période d'avril à juin (lignes 64 à 66) car la période d'estimation ne s'achève plus en janvier mais en mars.

III Table (26 fields, 66 records)							
違 <u>F</u> ile	📄 <u>E</u> dit 🛛 🕙	<u>G</u> enerate 🔣				0 X	
Table A	nnotations						
	\$TS-Market_4	\$TSLCI-Market_4	\$TSUCI-Market_4	\$TS-Total	\$TSLCI-Total	\$TSL	
47	13460.165	13046.567	13883.520	1895694.552	1890768.484	190 📥	
48	13637.234	13218.196	14066.159	1929821.249	1924806.501	193	
49	14038.478	13607.110	14480.023	1974007.314	1968877.747	197	
50	14588.176	14139.917	15047.010	2017063.960	2011822.507	202	
51	14826.444	14370.864	15292.773	2055709.852	2050367.976	206	
52	15328.900	14857.881	15811.032	2094273.974	2088831.887	209	
53	15403.883	14930.559	15888.373	2131431.902	2125893.258	213	
54	16187.796	15690.385	16696.942	2168729.836	2163094.271	217	
55	16303.304	15802.343	16816.083	2204919.579	2199189.973	221	
56	17250.576	16720.508	17793.149	2235223.381	2229415.030	224	
57	17616.290	17074.985	18170.366	2278910.104	2272988.230	228	
58	17639.270	17097.259	18194.069	2316079.288	2310060.827	232	
59	17552.150	17012.816	18104.209	2355228.381	2349108.190	236	
60	17499.120	16961.415	18049.510	2406836.211	2400581.914	241	
61	18183.056	17624.336	18754.958	2453038.341	2446663.985	245	
62	18512.777	17943.925	19095.050	2496354.087	2489867.172	250	
63	19125.395	18537.719	19726.936	2543477.283	2536867.916	255	
64	19394.782	18798.828	20004.796	2581510.338	2574802.140	258	
65	19387.631	18551.891	20251.298	2625230.895	2611195.788	263	
66	19550.898	18525.803	20617.962	2669744.972	2646565.409	269 👻	
	4		-				
						OK	

Figure 196. Table indiquant la nouvelle prévision

- Reliez un noeud de graphique Tracé horaire au nugget de modèle de séries temporelles.
 Cette fois, nous allons utiliser l'affichage Tracé horaire conçu spécifiquement pour les modèles de séries temporelles.
- 4. Dans l'onglet Tracé, affectez à Libellé de l'axe X la valeur Personnalisé, puis sélectionnez Date_.
- 5. Pour le Tracé, sélectionnez l'option Modèles de série temporelle sélectionnés.
- 6. Dans la liste **Série**, cliquez sur le bouton de sélection de champ, sélectionnez le champ \$TS-Market_4, puis cliquez sur **OK** pour l'ajouter à la liste.

📀 [Market_4 \$TS-Market_4 \$TSLCI-Market_4 \$TSUCI-Market_4] v. DATE_ 🛛 🔤
Plot Appearance Output Annotations
Plot: O Selected series 🖲 Selected Time Series models
Series:
X axis label: O Default O Custom
👿 Display series in separate panels 🛛 🐨 Normalize
Display: 📝 Line
Point
Smoother Smoother
🗹 Limit records Maximum number of records to plot: 2000 🗲
OK Run Cancel Apply Reset

Figure 197. Spécification des champs à tracer

7. Cliquez sur Exécuter.

Vous disposez maintenant d'un graphique qui illustre les ventes réelles pour Market_4 jusqu'à mars 2004, ainsi que les ventes prévisionnelles (Prévisions) et l'intervalle de confiance (indiqué par la zone ombrée bleue) jusqu'à juin 2004.

Comme dans le premier exemple, les valeurs prévisionnelles suivent étroitement les données réelles sur toute la période, ce qui indique une fois encore que vous disposez d'un modèle adéquat.



Figure 198. Prévision étendue jusqu'à juin

Récapitulatif

Vous avez appris à appliquer des modèles enregistrés afin d'étendre les prévisions antérieures lorsque de nouvelles données actuelles sont disponibles, et cela sans recréer vos modèles. Bien entendu, si vous avez une bonne raison de penser qu'un modèle a évolué, vous devez le recréer.

Chapitre 15. Prévision des ventes sur catalogue (séries temporelles)

Une société de vente sur catalogue souhaite prévoir les ventes mensuelles de sa ligne de vêtements pour hommes, en fonction des données des ventes réalisées au cours des 10 dernières années.

Cet exemple utilise le flux intitulé *catalog_forecast.str*, qui référence le fichier de données *catalog_seasfac.sav*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *catalog_forecast.str* se trouve dans le répertoire des *flux*.

Dans un exemple précédent, nous avons vu comment vous pouvez laisser Expert Modeler choisir à votre place le modèle le plus approprié pour les séries temporelles. L'heure est venue d'examiner de plus près les deux méthodes disponibles vous permettant de choisir vous-même un modèle, à savoir le lissage exponentiel et l'ARIMA.

Pour choisir un modèle approprié, il est d'abord recommandé de tracer les séries temporelles. L'inspection visuelle d'une série temporelle facilite souvent le choix à réaliser. Vous devez notamment vous poser les questions suivantes :

- La série présente-t-elle une tendance globale ? Si tel est le cas, la tendance est-elle constante ou s'atténue-t-elle avec le temps ?
- La série présente-t-elle des effets saisonniers ? Si tel est le cas, les fluctuations saisonnières croissent-elles avec le temps ou apparaissent-elles constantes sur des périodes successives ?

Création du flux

1. Créez un nouveau flux, puis ajoutez un noeud source Statistics pointant vers le fichier *catalog_seasfac.sav*.



Figure 199. Prévision des ventes sur catalogue

- 2. Ouvrez le noeud source IBM SPSS Statistics, puis cliquez sur l'onglet Types.
- 3. Cliquez sur Lire les valeurs, puis sur OK.
- 4. Cliquez sur la colonne Rôle du champ men, puis définissez le rôle sur Cible.

Catalog_seasfac.sav Image: Constraint of the seasfac.sav Image: CLEO_DEMOS/catalog_seasfac.sav Data Filter Types Annotations							
🔨 🕶 🗪 🕟 Read Values 🛛 Clear Values 🔹 Clear All Values							
Field -	Measurement	Values	Missing	Check	Role		
date	🔗 Continuous	[0000-12		None	O None	4	
🋞 men	🖉 Continuous	[3245.18,		None	O Target		
🛞 women	🔗 Continuous	[16578.9		None	O None		
🌮 jewel	🔗 Continuous	[5983.55,		None	O None		
🔿 mail	🔗 Continuous	[1147,15		None	O None		
🔆 page	🔗 Continuous	[51,114]		None	O None		
🗘 phone	🔗 Continuous	[17,59]		None	O None		
print	🔗 Continuous	[18061.2,		None	O None		
> service	🔗 Continuous	[15,68]		None	O None		
YEAR 🎆 Nominal 1989,199				None	O None		
Image: Second setting Image: Second setting							

Figure 200. Spécification du champ cible

- 5. Définissez le rôle de tous les autres champs sur Aucun, puis cliquez sur OK.
- 6. Reliez un noeud de graphique Tracé horaire au noeud source IBM SPSS Statistics.
- 7. Ouvrez le noeud Tracé horaire et, dans l'onglet Tracé, ajoutez men à la liste Série.
- 8. Affectez à Libellé de l'axe X la valeur Personnalisé, puis sélectionnez date.
- 9. Désélectionnez la case Normaliser.

😵 [men] v. date 📃 🔀
Plot Appearance Output Annotations
Plot: 💿 Selected series 🔘 Selected Time Series models
Series:
X axis label: 🔘 Default 🔘 Custom 🔗 date 🗐
📝 Display series in separate panels 🛛 Normalize
Display: 👿 Line
🥅 Point
🥅 Smoother
🗹 Limit records Maximum number of records to plot: 2000 🗧
OK Run Cancel Apply Reset

Figure 201. Traçage de la série temporelle

10. Cliquez sur Exécuter.

Analyse des données



Figure 202. Ventes réelles des vêtements pour hommes

La série indique une tendance ascendante générale ; en d'autres termes, les valeurs de la série ont tendance à augmenter dans le temps. La tendance ascendante semble constante, ce qui indique une tendance linéaire.

En outre, la série présente un schéma saisonnier caractérisé par des ventes élevées en décembre, comme l'indiquent les lignes verticales du graphique. Les variations saisonnières semblent croître avec la tendance ascendante de la série, ce qui suggère la présence d'effets saisonniers multiplicatifs plutôt qu'additifs.

1. Cliquez sur **OK** pour fermer le graphique.

Les caractéristiques de la série étant identifiées, vous pouvez essayer de la modéliser. La méthode de lissage exponentiel permet de prévoir les séries qui présentent une tendance et/ou des effets saisonniers. Comme nous l'avons vu, les données présentent les deux caractéristiques.

Lissage exponentiel

La création d'un modèle de lissage exponentiel offrant le meilleur ajustement implique la détermination du type de modèle (à savoir, si le modèle doit inclure la tendance et/ou les effets saisonniers), puis l'obtention des paramètres offrant le meilleur ajustement pour le modèle choisi.

Le graphique des ventes de vêtements pour hommes dans le temps suggérait un modèle avec à la fois une composante de tendance linéaire et une composante d'effets saisonniers multiplicatifs. Cela implique un modèle de Winters. Toutefois, dans un premier temps, nous explorerons un modèle simple (sans tendance, ni effets saisonniers), puis un modèle de Holt (incorporant une tendance linéaire, mais aucun effet saisonnier). Au terme de cette opération, vous saurez identifier un modèle qui ne constitue pas un ajustement aux données adéquat, compétence essentielle pour la création réussie de modèles.

Fields	Data Speci	ifications	Build Options	Model Options	Annotations
<u>S</u> elect ar	n item:	-			
Genera	Ľ.	Method	Exponential S	Smoothing 👻	
Output		- Model ◎ S ◎ H ◎ B ◎ D	Type imple folt's linear tren rown's linear tre amped trend	O Simple s d O Winters' end O Winters'	easonal additive multiplicative

Figure 203. Spécification du lissage exponentiel

Nous allons commencer avec un modèle de lissage exponentiel simple.

- 1. Ajoutez un noeud Séries temporelles au flux et reliez-le au noeud source.
- 2. Dans l'onglet Spécifications des données de la sous-fenêtre Observations, sélectionnez date comme champ Date/Heure.
- 3. Sélectionnez Mois comme Intervalle de temps.

📀 men					
Fields	Data Specifications	Build Options	Model Options	Annotations	
<u>S</u> elect ar	n item:				
Observa	ations	🔵 💿 <u>O</u> bserva	ations are specifi	ed by a date/tir	me field
Time Inf	terval	Date/tir	me field:		
Aggrega	ation and Distribution	🔗 dat	e		-1
Missing	Value Handling	T <u>i</u> me ir	nterval: Months		*

Figure 204. Définition de l'intervalle de temps

- 4. Dans l'onglet Options de création de la sous-fenêtre Général, définissez l'option Méthode sur Lissage exponentiel.
- 5. Définissez Type de modèle sur Simple.

Fields	Data Speci	ifications	Build Options	Model Options	Annotations
<u>S</u> elect ar	n item:				
General	ť (Method	Exponential S	Smoothing 👻	
Output		- Model ◎ S ◎ F ◎ B ◎ C	Type imple lolt's linear tren rown's linear tre vamped trend	◯ Simple s d ◯ Winters' a end ◯ Winters' r	easonal additive multiplicative

Figure 205. Définition de la méthode de génération de modèle

6. Cliquez sur Exécuter pour créer le nugget de modèle.

📀 [men \$TS-men] v. date 🧱 🛃
Plot Appearance Output Annotations
Plot: 💿 Selected series 🔘 Selected Time Series models
Series: Series:
X axis label: 🔘 Default 🔘 Custom 🛷 date 🗐
🥅 Display series in separate panels 🛛 🕅 Normalize
Display: 👿 Line
🥅 Point
🧾 Smoother
🗹 Limit records Maximum number of records to plot: 2000 🤤
OK Run Cancel Apply Reset

Figure 206. Traçage du modèle de séries temporelles

- 7. Reliez un noeud Tracé horaire au nugget de modèle.
- 8. Dans l'onglet Tracé, ajoutez men et \$TS-men à la liste Série.
- 9. Affectez à Libellé de l'axe X la valeur Personnalisé, puis sélectionnez date.

10. Désélectionnez les cases Afficher les séries dans des panneaux distincts et Normaliser.

40000 30000 20000 10000 1988-01-01 1990-01-01 1994-01-01 1996-01-01 1998-01-01 2000-01-01 date

11. Cliquez sur **Exécuter**.

Figure 207. Modèle de lissage exponentiel simple

Le tracé **men** représente les données réelles, tandis que **\$TS-men** désigne le modèle de séries temporelles.

Bien que le modèle simple présente, en fait, une tendance ascendante progressive (et plutôt lourde), il ne prend aucunement en compte les effets saisonniers. Vous pouvez exclure ce modèle en toute sécurité.

12. Cliquez sur OK pour fermer la fenêtre du tracé horaire.

Method: Exponential Smo	pothing 🔫
Model Type	
O Simple	🔘 Simple seasonal
Holt's linear trend	🔘 Winters' additive
O Brown's linear trend	O Winters' multiplicative
O Damped trend	

Figure 208. Sélection du modèle de Holt

Essayons le modèle linéaire de Holt. Celui-ci devrait au moins mieux modéliser la tendance que le modèle simple, bien que lui aussi ne soit probablement pas en mesure de capturer les effets saisonniers.

- 13. Rouvrez le noeud Séries temporelles.
- 14. Dans l'onglet Options de création de la sous-fenêtre Général, Lissage exponentiel étant toujours sélectionné comme Méthode, sélectionnez Tendance linéaire de Holt comme Type de modèle.
- 15. Cliquez sur Exécuter pour recréer le nugget de modèle.
- 16. Réouvrez le noeud Tracé horaire et cliquez sur Exécuter.



Figure 209. Modèle de tendance linéaire de Holt

Le modèle de Holt affiche une tendance ascendante plus lisse que le modèle simple, mais ne prend aucunement en compte les effets saisonniers ; par conséquent, vous pouvez également le supprimer.

17. Fermez la fenêtre du tracé horaire.

Vous vous souvenez peut-être que le graphique initial des ventes de vêtements pour hommes dans le temps suggérait un modèle incorporant une tendance linéaire et des effets saisonniers multiplicatifs. Par conséquent, le modèle de Winters pourrait être plus approprié.

Fields	Data Speci	fications	Build Options	Model Options	Annotations
<u>S</u> elect ar	n item:				
General	ř	Method	Exponential S	Smoothing 🔫	
Output		- Model ◎ S ◎ F ◎ B ◎ D	Type iimple folt's linear tren rown's linear tre pamped trend	◯ Simple s d ◯ Winters' a end ⓒ Winters' r	easonal additive nultiplicative

Figure 210. Sélection du modèle de Winters

- 18. Rouvrez le noeud Séries temporelles.
- 19. Dans l'onglet Options de création de la sous-fenêtre Général, Lissage exponentiel étant toujours sélectionné comme Méthode, sélectionnez Méthode multiplicative de Winters comme Type de modèle.
- 20. Cliquez sur Exécuter pour recréer le nugget de modèle.
- 21. Ouvrez le noeud Tracé horaire et cliquez sur Exécuter.



Figure 211. Modèle multiplicatif de Winters

Ce modèle semble plus approprié car il reflète à la fois la tendance et les effets saisonniers des données.

Le jeu de données couvre une période de 10 années et comprend 10 pics saisonniers, se produisant au mois de décembre de chaque année. Les 10 pics présents dans les résultats prévus correspondent parfaitement aux 10 pics annuels dont font état les données réelles.

Toutefois, les résultats soulignent les limites de la procédure de lissage exponentiel. L'observation des pointes ascendantes et descendantes révèle l'existence d'une structure significative non expliquée.

Si vous êtes essentiellement intéressé par la modélisation d'une tendance à long terme avec variation saisonnière, le lissage exponentiel peut s'avérer un choix adéquat. Pour modéliser une structure plus complexe telle que celle-ci, nous devons envisager l'utilisation de la procédure ARIMA.

ARIMA

Grâce à la procédure ARIMA, vous pouvez créer un modèle ARIMA (AutoRegressive Integrated Moving Average) qui vous permet d'affiner la modélisation des séries temporelles. Les modèles ARIMA mettent à votre disposition des méthodes de modélisation des composantes de tendance et saisonnières plus élaborées que celles proposées par les modèles de lissage exponentiel, et ils permettent d'inclure des variables de prédicteur dans le modèle.

A travers l'exemple de la société de vente sur catalogue qui souhaite développer un modèle de prévision, nous avons vu comment celle-ci a collecté des données sur les ventes mensuelles de vêtements pour hommes, ainsi que plusieurs séries susceptibles de faciliter l'explication d'une partie de la variation des ventes. Parmi les prédicteurs possibles figurent le nombre de catalogues envoyés par messagerie électronique, le nombre de pages dans le catalogue, le nombre de lignes téléphoniques dédiées à la prise de commandes, les frais de publicité imprimée et le nombre de conseillers du service à la clientèle.

Certains de ces prédicteurs sont-ils utiles pour la prévision ? Un modèle comportant des prédicteurs est-il réellement meilleur qu'un modèle qui en est dépourvu ? A l'aide de la procédure ARIMA, nous pouvons créer un modèle de prévision comportant des prédicteurs et déterminer s'il existe une différence significative en matière de capacité prédictive par rapport au modèle de lissage exponentiel dépourvu de prédicteur.

Avec la méthode ARIMA, vous pouvez affiner le modèle en spécifiant les ordres d'autorégression, de différenciation et de moyenne mobile, ainsi que les équivalents saisonniers de ces composantes. Etant

donné que la détermination manuelle des meilleures valeurs pour ces composantes peut être un processus de longue durée impliquant une série d'essais et d'erreurs, nous allons, pour cet exemple, laisser Expert Modeler choisir un modèle ARIMA à notre place.

Nous allons essayer de créer un meilleur modèle en traitant certaines des autres variables du jeu de données en tant que variables de prédicteur. Celles qui semblent le plus utile à inclure en tant que prédicteurs sont le nombre de catalogues envoyés par messagerie électronique (mail), le nombre de pages dans le catalogue (page), le nombre de lignes téléphoniques dédiées à la prise de commandes (phone), les frais de publicité imprimée (print) et le nombre de conseillers du service à la clientèle (service).

Catalog_seasfac.sav Preview Refresh \$CLEO_DEMOS/catalog_seasfac.sav							
Data Filter Types Annotations Image: Second							
Field -	Measurement	Values	Missing	Check	Role		
date	🔗 Continuous	[0000-12	-	None	○ None	4	
🛞 men	🔗 Continuous	[3245.18,		None	Target		
🛞 women	🔗 Continuous	[16578.9		None	O None		
🛞 jewel	🔗 Continuous	[5983.55,	1	None	O None		
🔷 mail	🔗 Continuous	[1147,15		None	🔪 Input		
🔆 page	🔗 Continuous	[51,114]		None	🔪 Input		
🔆 phone	🔗 Continuous	[17,59]		None	🔪 Input		
🛞 print	🔗 Continuous	[18061.2,		None	🔪 Input		
🔆 service	🔗 Continuous	[15,68]		None	🔪 Input		
🔆 YEAR_	💑 Nominal	1989,199		None	🛇 None	-	
Image: Second setting							

Figure 212. Définition des champs de prédicteur

- 1. Ouvrez le noeud source de fichier IBM SPSS Statistics.
- Dans l'onglet Types, définissez l'option Rôle des variables mail, page, phone, print et service sur Entrée.
- **3**. Vérifiez que la direction pour men est définie sur **Cible** et que tous les champs restants sont définis sur **Aucun**.
- 4. Cliquez sur OK.
- 5. Ouvrez le noeud Séries temporelles.

- 6. Dans l'onglet Options de création de la sous-fenêtre Général, définissez l'option Méthode sur Expert Modeler.
- 7. Sélectionnez l'option Modèles ARIMA uniquement, puis vérifiez que la case Expert Modeler prend en compte les modèles saisonniers est cochée.

Fields Data Spe	cifications	Build Options	Model Options	Annotations			
<u>S</u> elect an item:	_						
General	General Method: Expert Modeler 🔫						
Output	Model A E A A M A A A	Type II models xponential smo RIMA models o Expert Modeler (oothing models o nly considers seaso	inly nal models			

Figure 213. Sélection des modèles ARIMA uniquement

- 8. Cliquez sur Exécuter pour recréer le nugget de modèle.
- 9. Ouvrez le nugget de modèle.

Dans la colonne de gauche de l'onglet Sortie, sélectionnez les **Informations sur le modèle**. Comme vous pouvez le constater, Expert Modeler n'a choisi, parmi les cinq prédicteurs spécifiés, que deux valeurs comme étant significatives pour le modèle.



Figure 214. Expert Modeler choisit deux prédicteurs

- 10. Cliquez sur OK pour fermer le nugget du modèle.
- 11. Ouvrez le noeud Tracé horaire et cliquez sur **Exécuter**.



Figure 215. Modèle ARIMA avec prédicteurs spécifiés

Ce modèle est meilleur que le précédent car il capture également la grande pointe descendante, ce qui en fait le meilleur ajustement constaté jusqu'à présent.

Nous pourrions essayer d'affiner davantage le modèle, mais toute amélioration à partir de ce point est susceptible d'être minimale. Nous avons démontré que le modèle ARIMA avec prédicteurs est préférable ; par conséquent, utilisons le modèle que nous venons de créer. Dans le cadre de cet exemple, nous allons prévoir les ventes de l'année à venir.

- 12. Cliquez sur OK pour fermer la fenêtre du tracé horaire.
- 13. Ouvrez le noeud Séries temporelles et sélectionnez l'onglet Options de modèle.
- 14. Cochez la case Etendre les enregistrements dans le futur, puis attribuez-lui la valeur 12.
- 15. Cochez la case Calculer les valeurs futures des entrées.
- 16. Cliquez sur Exécuter pour recréer le nugget de modèle.
- 17. Ouvrez le noeud Tracé horaire et cliquez sur Exécuter.

La prévision pour 1999 semble correcte : comme prévu, il existe un retour aux niveaux normaux de ventes après le pic de décembre, ainsi qu'une tendance ascendante régulière dans la seconde moitié de l'année, dont les ventes sont globalement supérieures à celles de l'année précédente.



Figure 216. Prévision des ventes prolongée de 12 mois

Récapitulatif

Vous avez correctement modélisé une série temporelle complexe, en intégrant non seulement une tendance ascendante, mais aussi des effets saisonniers et d'autres variations. Vous avez également vu comment, à l'aide d'une série d'essais et d'erreurs, vous avez pu approcher petit à petit un modèle précis, que vous avez ensuite utilisé pour prévoir les ventes à venir.

Dans la pratique, vous devriez réappliquer le modèle car les données de ventes réelles sont mises à jour, par exemple, chaque mois ou trimestre, et générer des prévisions actualisées. Pour plus d'informations, voir la rubrique «Réapplication d'un modèle de séries temporelles», à la page 171.

Chapitre 16. Propositions aux clients (auto-apprentissage)

Le noeud de réponse Auto-formation permet de générer et de mettre à jour un modèle grâce auquel vous pouvez prévoir les offres les plus appropriées pour les clients et la probabilité d'acceptation des offres. Ces types de modèle sont les plus utiles dans la gestion de la relation client, notamment dans les applications marketing ou les centres d'appels.

Cet exemple est basé sur un établissement bancaire fictif. Le service marketing souhaite obtenir des résultats plus rentables lors des prochaines campagnes en présentant à chaque client une offre de service financier adaptée. Plus précisément, l'exemple utilise un modèle de réponse d'auto-apprentissage pour identifier les caractéristiques des clients les plus susceptibles de répondre favorablement, en fonction d'offres et de réponses précédentes, et de promouvoir la meilleure offre existant à partir des résultats.

Cet exemple utilise le flux *pm_selflearn.str* qui fait référence aux fichiers de données *pm_customer_train1.sav, pm_customer_train2.sav* et *pm_customer_train3.sav*. Ces fichiers sont disponibles dans le dossier *Demos* de toutes les installations d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *pm_selflearn.str* se trouve dans le dossier des *flux*.

Données existantes

L'entreprise dispose de données historiques retraçant les offres faites aux clients lors des campagnes précédentes, ainsi que les réponses à ces offres. Ces données incluent également des informations démographiques et financières qui peuvent être utilisées pour prévoir les taux de réponse pour différents clients.

違 File	📄 <u>E</u> dit 🛛 💐) <u>G</u> enerate						0
Table ,	Annotations							
	customer_id	campaign	response	response_date	purchase	purchase_date	product_id	Rowid
1	7	2	0	\$null\$	0	\$null\$	\$null\$	1
2	13	2	0	\$null\$	0	\$null\$	\$null\$	2
3	15	2	0	\$null\$	0	\$null\$	\$null\$	3
4	16	2	1	2006-07-05 00:00:00	0	\$null\$	183	761
5	23	2	0	\$null\$	0	\$null\$	\$null\$	4
6	24	2	0	\$null\$	0	\$null\$	\$null\$	5
7	30	2	0	\$null\$	0	\$null\$	\$null\$	6
8	30	3	0	\$null\$	0	\$null\$	\$null\$	7
9	33	2	0	\$null\$	0	\$null\$	\$null\$	8
10	42	3	0	\$null\$	0	\$null\$	\$null\$	9
11	42	2	0	\$null\$	0	\$null\$	\$null\$	10
12	52	2	0	\$null\$	0	\$null\$	\$null\$	11
13	57	2	0	\$null\$	0	\$null\$	\$null\$	12
14	63	2	1	2006-07-14 00:00:00	0	\$null\$	183	1501
15	74	2	0	\$null\$	0	\$null\$	\$null\$	13
16	74	3	0	\$null\$	0	\$null\$	\$null\$	14
17	75	2	0	\$null\$	0	\$null\$	\$null\$	15
18	82	2	0	\$null\$	0	\$null\$	\$null\$	16
19	89	3	0	\$null\$	0	\$null\$	\$null\$	17
20	89	2	0	\$null\$	0	\$null\$	\$null\$	18
						too hard a second		1

Figure 217. Réponses aux offres précédentes

Création du flux

1. Ajoutez un noeud source de fichier Statistics qui pointe sur *pm_customer_train1.sav*, dans le dossier *Demos* du répertoire d'installation d'IBM SPSS Modeler.



Figure 218. Flux d'échantillons MRAA

- 2. Ajoutez un noeud Remplacer et saisissez campaign comme champ à renseigner.
- 3. Sélectionnez Toujours comme type de champ Remplacer.
- 4. Dans la zone de texte Remplacer par, entrez to_string(campaign) et cliquez sur OK.

🚱 Filler	X
Preview	
Settings Annotations	
Fill in fields:	
🖋 campaign	
Replace: Always Condition:	
@BLANK(@FIELD)	
Replace with:	
to_string(campaign)	
OK Cancel	Apply Reset

Figure 219. Calcul d'un champ de campagne

5. Ajoutez un noeud type et définissez l'option *Rôle* sur **Aucun** pour les champs *customer_id*, *response_date*, *purchase_date*, *product_id*, *Rowid*, et *X_random*.

Type	iew		_	(0
~	Read Values	Clear V	alues	Clear All Valu	Jes
Field -	Measurement	Values	Missing	Check	Role
🔆 customer_id	🖉 Continuous	[7,116993]		None	🛇 None 🤘
A campaign	Nominal	"1","2","3		None	O Target
🔆 response	🎖 Flag	1/0		None	O Target
🔁 response_date	🔗 Continuous	[2006-04		None	O None
🔆 purchase	💑 Nominal	0,1		None	🔪 Input
🔁 purchase_date	🔗 Continuous	[2006-04		None	○ None
product_id	🖉 Continuous	[183,421]		None	○ None
🔆 Rowid	🖉 Continuous	[1,19599]		None	○ None
🔆 age	🖉 Continuous	[10,96]		None	🔪 Input
🔉 age_youngest	🖉 Continuous	[0,66]		None	🔪 Input 🗖
View current fiel OK Cancel	ds 🔘 View unused	field settings		ĺ	Apply Reset

Figure 220. Modification des paramètres du noeud type

6. Définissez l'option *Rôle* sur **Cible** pour les champs *campaign* et *response*. Il s'agit des champs sur lesquels doivent reposer les prévisions.

Définissez la mesure sur Indicateur pour le champ response.

7. Cliquez sur Lire les valeurs, puis sur OK.

Comme les données de champ de campagne se présentent sous la forme d'une liste numérique (1, 2, 3 et 4), vous pouvez recodifier les champs pour obtenir des titres plus évocateurs.

- 8. Ajoutez un noeud Recoder au noeud type.
- 9. Dans le champ Recoder dans, sélectionnez Champ existant.
- 10. Dans la liste Recoder le champ, sélectionnez campagne.
- 11. Cliquez sur le bouton Obtenir ; les valeurs de campagne sont ajoutées à la colonne Valeur d'origine.
- **12**. Dans la colonne *Nouvelle valeur*, entrez les noms de campagne suivants sur les quatre premières lignes :
 - Prêt hypothécaire
 - Prêt automobile
 - Epargne
 - Retraite
- 13. Cliquez sur OK.

Reclassify	review		
**			
Settings Annot	ations		
	Mode:	💿 Single 🔘 Multiple	
	Reclassify into:	🔘 New field 💿 Existing field	E.
Reclassify field:			
💑 campaign			-
New field name:			
Reclassify3			
Reclassify value:	5:		
🕨 🕨 Get	💛 Copy	🧷 Clear new	🗳 Auto
Origin	al value 🗂	New value	
1		Mortgage	
2		Car Ioan	
3		Savings	
4		Pension	-
For unspecified v	values use: 🔘 O	riginal value 🔘 Default value	undef
04			(Inch) Dec

Figure 221. Recodification des champs de campagne

14. Reliez un noeud de modélisation MRAA au noeud Recoder. Dans l'onglet Champs, sélectionnez **campaign** pour le champ Cible et **response** pour le champ de réponse Cible.

😡 campaign			X
Fields Model Settings	Annotations		
Target field:	o campaign	ا ب]	
Target response field:	response		
I Use type node settings	O Use custom	settings	
Inputs:		-	
		×	
_			
Partition:		▼ 1	1
Use frequency field		-	1
OK 🕨 Run Cano	el	Apply	t

Figure 222. Sélection de la cible et de la réponse cible

15. Dans le champ Nombre maximal de prédictions par enregistrement de l'onglet Paramètres, réduisez la valeur à 2.

Deux offres seront ainsi identifiées comme présentant la plus forte probabilité d'acceptation pour chaque client.

16. Assurez-vous de sélectionner Prendre en compte la fiabilité du modèle et cliquez sur Exécuter.

😡 campaign			X
		0	-0
Fields Model Settings Annota	tions		
Maximum number of predictions per re	ecord: 2		
Level of randomization :	0.00 ≑		
Set random seed:	876547 🚔		
Sort order:			
Descending(offers with hig	hest score will be returned	Ð	
O Ascending(offers with low	est score will be returned)		
Preferences for target fields:	10		
Value	Preference	Always include	Add
			Delete
Vake account of model reliability			
OK 🕨 Run Cancel		Apply	Reset

Figure 223. Paramètres du noeud MRAA

Navigation dans le modèle

1. Ouvrez le nugget de modèle. Initialement, l'onglet Modèle indique l'estimation de l'exactitude des prévisions pour chaque offre, ainsi que l'importance relative de chaque prédicteur pour estimer le modèle.

Pour afficher la corrélation de chaque prédicteur avec la variable cible, choisissez **Association avec réponse** dans la liste **Vue** du panneau de droite.

2. Pour vous déplacer entre les quatre offres auxquelles les prédictions s'appliquent, sélectionnez l'offre appropriée dans la liste **Vue** du panneau de gauche.



Figure 224. Nugget de modèle MRAA

- 3. Fermez la fenêtre du nugget de modèle.
- 4. Dans l'espace de travail, déconnectez le noeud source IBM SPSS Statistics pointant vers le fichier *pm_customer_train1.sav*.
- 5. Ajoutez un noeud source de fichier Statistics pointant vers le fichier *pm_customer_train2.sav* situé dans le dossier *Demos* du répertoire d'installation d'IBM SPSS Modeler, puis connectez-le au noeud Remplacer.



pm_customer_train3.s..

Figure 225. Association d'une deuxième source de données au flux MRAA

6. Dans l'onglet Modèle du noeud MRAA, sélectionnez Poursuivre le modèle d'apprentissage existant.

💟 campaign	\mathbf{X}
Fields Model Settings Annotations	
Model name: O Auto Custom	
Use partitioned data	
Continue training existing model	
Target field values: () Use all () Specify	
	Add
	Edit
	Delete
Model Assessment	
☑ Include model assessment	
Set random seed: 876547 🗲	
Simulated sample size:	
Display model evaluation	
OK Run Cancel Apply	Reset

Figure 226. Poursuite de l'apprentissage du modèle

7. Cliquez sur **Exécuter** pour recréer le nugget de modèle. Pour en afficher les détails, double-cliquez sur le nugget dans l'espace de travail.

L'onglet Modèle affiche les estimations révisées de l'exactitude des prévisions pour chaque offre.

8. Ajoutez un noeud source de fichier Statistics pointant vers le fichier *pm_customer_train3.sav* situé dans le dossier *Demos* du répertoire d'installation d'IBM SPSS Modeler, puis connectez-le au noeud Remplacer.



pm_customer_train3.s..

Figure 227. Association d'une troisième source de données au flux MRAA

- **9**. Cliquez sur **Exécuter** pour recréer une nouvelle fois le nugget de modèle. Pour en afficher les détails, double-cliquez sur le nugget dans l'espace de travail.
- **10**. L'onglet Modèle affiche maintenant les estimations finales de l'exactitude des prévisions pour chaque offre.

Comme vous pouvez le constater, l'exactitude moyenne a légèrement baissé (de 86,9% à 85,4%) à la suite de l'ajout des sources de données supplémentaires. Cela étant, cette faible variation peut s'expliquer par la présence de petites anomalies dans les données disponibles.



Figure 228. Nugget de modèle MRAA mis à jour

- 11. Reliez un noeud Table au dernier (troisième) modèle généré, puis exécutez ce noeud.
- **12**. Faites défiler le tableau vers la droite. Les prévisions indiquent les offres qu'un client est le plus enclin à accepter et la confiance qu'il sera enclin à vous accorder en fonction des détails le concernant.

Par exemple, selon la première ligne du tableau illustré, la confiance avec laquelle un client ayant déjà contracté un prêt automobile acceptera une éventuelle offre d'épargne retraite s'élèvera à 13,2 % (indiqué par la valeur 0,132 dans la colonne *\$SC-campaign-1*). Toutefois, les deuxième et troisième lignes présentent deux autres clients ayant souscrit ce même type de prêt ; dans leur cas, il existe un niveau de fiabilité de 95,7 % qu'ils ouvriraient un compte d'épargne suite à une offre, et un niveau de fiabilité à 80% qu'ils accepteraient une offre d'épargne de retraite.

💫 File 📄 Edit 🕙 Generate 🔣 🕒 📢 🛗							
Table Annotations							
		X_random	\$S-campaign-1	\$SC-campaign-1	\$S-campaign-2	\$SC-campaign-2	
1//////////////////////////////////////		1	Pension	0.132	Mortgage	0.107	
2		1	Savings	0.957	Pension	0.844	
3		1	Savings	0.957	Pension	0.802	
4		3	Pension	0.132	Mortgage	0.107	
5		1	Pension	0.805	Savings	0.284	
6		3	Pension	0.132	Mortgage	0.107	
7		2	Pension	0.132	Mortgage	0.107	
8		3	Pension	0.132	Mortgage	0.107	
9		1	Pension	0.132	Mortgage	0.107	
10		1	Pension	0.132	Mortgage	0.107	
11		2	Pension	0.132	Mortgage	0.107	
12		2	Pension	0.132	Mortgage	0.107	
13		2	Savings	0.957	Mortgage	0.829	
14		2	Savings	0.164	Pension	0.132	
15		2	Savings	0.957	Pension	0.868	
16		2	Pension	0.132	Mortgage	0.107	
17		3	Pension	0.132	Mortgage	0.107	
18		3	Pension	0.132	Mortgage	0.107	
19		3	Savings	0.289	Pension	0.132	
20		2	Pension	0.132	Mortgage	0.107	
	1						

Figure 229. Résultat du modèle : prédictions d'offre et de confiance

Des explications sur les fondements mathématiques des méthodes de modélisation utilisées dans IBM SPSS Modeler sont présentées dans le *guide des algorithmes d'IBM SPSS Modeler*, disponible sous forme de fichier PDF avec le téléchargement de votre produit.

Sachez également que ces résultats sont établis uniquement d'après les données d'apprentissage. Pour évaluer comment le modèle peut se généraliser à d'autres données dans le monde réel, vous devez utiliser un noeud Partitionner destiné à contenir un sous-ensemble d'enregistrements à des fins de test et de validation.

Chapitre 17. Prévision des défauts de paiement (Réseau Bayésien)

Les réseaux Bayésiens permettent de créer un modèle de probabilité en combinant les preuves observées et enregistrées avec les connaissances réelles "de bon sens" pour établir la probabilité des occurences en utilisant des attributs apparemment sans lien.

Cet exemple utilise le flux nommé *bayes_bankloan.str*, qui fait référence au fichier de données *bankloan.sav*. Ces fichiers sont accessibles dans le répertoire *Demos* de toute installation IBM SPSS Modeler. Vous pouvez y accéder à partir du groupe de programmes IBM SPSS Modeler du menu Démarrer de Windows. Le fichier *bayes_bankloan.str* se trouve dans le répertoire des *flux*.

Par exemple, imaginons qu'une banque s'inquiète des prêts susceptibles de ne pas être remboursés. Si les données par défaut des prêts précédents peuvent être utilisées pour prédire quels clients potentiels risquent d'avoir des difficultés à rembourser leur prêt, il est alors possible de refuser un prêt à ces clients à "fort risque" ou de leur proposer d'autres produits.

Cet exemple se concentre sur l'utilisation de données par défaut sur des prêts existants permettant de prédire quels futurs clients sont susceptibles de ne pas pouvoir rembourser leur prêt et examine trois différents types de modèle de réseau Bayésien afin d'établir celui le plus adapté à cette situation.

Création du flux

1. Ajoutez un noeud source Statistics pointant sur bankloan.sav dans le dossier Demos.



Figure 230. Flux d'échantillons de réseau Bayésien

- 2. Ajoutez un noeud type au noeud source et définissez le rôle du champ par **défaut** sur **Cible**. Le rôle de tous les autres champs doit être défini sur **Entrée**.
- 3. Cliquez sur le bouton Lire les valeurs pour remplir la colonne Valeurs.

Types Forma	view t Annotations			(
۰.	🗪 🕩 Read Va	alues Cle	ear Values	Clear A	I Values
Field	Measurement	Values	Missing	Check	Role
🗘 age	Continuous	[20,56]		None	🔪 Input
🔆 ed	Ordinal	1,2,3,4,5		None	> Input
🔆 employ	Continuous	[0,33]		None	🔪 Input
🗘 address	Continuous	[0,34]		None	🔪 Input
income	Continuous	[13.0,44		None	🔪 Input
debtinc	Continuous	[0.1,41.3]		None	🔪 Input
creddebt	Continuous	[0.01169		None	🔪 Input
othdebt	Continuous	[0.04558		None	🔪 Input
🔆 default	🎖 Flag	1/0		None	Target
OK Can	t fields 🔘 View uni	used field se	ttings	Ap	ply <u>R</u> ese

Figure 231. Sélection du champ cible

Les observations ayant une valeur de cible nulle ne servent à rien lors de la création d'un modèle. Il est possible d'exclure ces observations afin qu'elles ne soient pas utilisées lors de l'évaluation du modèle.

- 4. Ajoutez un noeud Sélectionner au noeud type.
- 5. Pour le mode, choisissez **Supprimer**.
- 6. Dans la boîte de dialogue Condition, saisissez valeur par défault = '\$null\$'.

Select	
	0
Settings Annotations	
Mode: 🔘 Include 💿 Discard	
default = '\$null\$' Condition:	
OK Cancel	Apply Reset

Figure 232. Suppression des cibles à valeur nulle

Il est possible de créer plusieurs types de réseaux Bayésiens. Par conséquent, il peut être utile d'en comparer plusieurs afin de connaître celui qui offre les meilleures prédictions. Le premier est un modèle Tree Augmented Naïve Bayes (TAN).

7. Reliez un noeud de réseau Bayésien au noeud Sélectionner.

- 8. Pour choisir le nom du modèle, dans l'onglet Modèle, sélectionnez **Personnalisé** puis saisissez TAN dans la zone de texte.
- 9. Pour le type de structure, sélectionnez TAN et cliquez sur OK.

TAN			X
			0.0
<u> </u>			
Fields Model Exper	t Analyze Annotations		
Model name:	🔘 Auto 🔘 Custom	TAN	
🔽 Use partitioned data			
👿 Build model for each	n split		
To select fields manual	ly, choose "Use custom settin	ngs" on the Fields tab —	
Partition:			-
Splits:			×
Continue training exis	ting model	-1-+	
Structure type:		nket	
Barameter learning meth	tion preprocessing step		
Farameter lean ling meth	🔘 Maximum likelihood 🔇	Bayes adjustment for	small cell counts
OK 🕨 Run Ca	ncel		Apply Reset

Figure 233. Création d'un modèle Tree Augmented Naïve Bayes

Le deuxième type de modèle à créer a une structure de couverture de Markov.

- 10. Relier un deuxième noeud de réseau Bayésien au noeud Sélectionner.
- 11. Pour choisir le nom du modèle, dans l'onglet Modèle, sélectionnez **Personnalisé** puis saisissez Markov dans la zone de texte.
- 12. Pour le type de structure, sélectionnez Couverture de Markov et cliquez sur OK.

🚰 Markov	
	0
Fields Model Expert Analyze Annotations	
Model name: O Auto O Custom Markov	
☑ Use partitioned data	
☑ Build model for each split	
To select fields manually, choose "Use custom settings" on the Fields ta	b
Partition:	-
Splits:	×
Continue training existing model	
Structure type: O TAN I Markov Blanket	
Include feature selection preprocessing step	
Parameter learning method:	
🔘 Maximum likelihood 🔘 Bayes adjustmen	t for small cell counts
OK FRUN Cancel	Apply Reset

Figure 234. Création d'un modèle Couverture de Markov

Le troisième type de modèle a une structure de couverture de Markov et utilise également le prétraitement de la sélection de fonctions pour sélectionner les entrées importantes liées à la variable cible.

- 13. Relier un troisième noeud de réseau Bayésien au noeud Sélectionner.
- 14. Pour choisir le nom du modèle, dans l'onglet Modèle, sélectionnez **Personnalisé** puis saisissez Markov-FS dans la zone de texte.
- 15. Pour le type de structure, sélectionnez Couverture de Markov.
- 16. Sélectionner Inclure une étape de prétraitement de la sélection des fonctions et cliquer sur OK.
| 💟 Markov-FS | $\overline{\mathbf{X}}$ |
|---|-------------------------|
| | |
| Fields Model Expert Analyze Annotations | |
| Model name: O Auto O Custom Markov-FS | |
| 👿 Use partitioned data | |
| ☑ Build model for each split | |
| _To select fields manually, choose "Use custom settings" on the Fields tab- | |
| Partition: | - |
| Splits: | × |
| Continue training existing model | |
| Structure type: 🔘 TAN 🔘 Markov Blanket | |
| ☑ Include feature selection preprocessing step | |
| Parameter learning method: | |
| ◙ Maximum likelihood ◎ Bayes adjustment f | or small cell counts |
| OK Run Cancel | Apply Reset |

Figure 235. Création d'un modèle de couverture de Markov avec prétraitement de la sélection des fonctions

Navigation dans le modèle

1. Exécutez le flux pour créer des nuggets de modèle, qui sont ajoutés au flux et à la palette Modèles dans l'angle supérieur droit. Pour afficher leurs détails, double-cliquez sur l'un des nuggets de modèle du flux.

L'onglet Modèle du nugget de modèle est divisé en deux panneaux. Le panneau de gauche contient un graphique de noeuds en réseau qui affiche la relation entre la cible et ses prédicteurs les plus importants, ainsi que la relation entre les prédicteurs.

Le panneau de droite affiche soit l'*Importance des prédicteurs*, qui indique l'importance relative de chaque prédicteur pour l'estimation du modèle ou les *Probabilités conditionnelles*, qui contiennent la valeur de la probabilité conditionnelle de chaque valeur de noeud et pour chaque combinaison de valeurs dans ses noeuds parent.



Figure 236. Affichage d'un modèle Tree Augmented Naïve Bayes

- 2. Connectez le nugget de modèle TAN au nugget Markov (choisissez **Remplacer** dans la boîte de dialogue d'avertissement).
- **3.** Connectez le nugget Markov au nugget Markov-FS (choisissez **Remplacer** dans la boîte de dialogue d'avertissement).
- 4. Alignez les trois nuggets avec le noeud Sélectionner pour une meilleure lecture.



Figure 237. Alignement des nuggets dans le flux

- 5. Pour renommer les sorties de modèle et obtenir un graphique d'évaluation plus clair, liez le noeud Filtrer au nugget de modèle Markov-FS.
- 6. Dans la colonne de droite *Champ*, remplacez le nom \$B par défaut par TAN, le nom \$B1 par défaut par Markov et le nom \$B2 par défaut par Markov-FS.

Filter		
7	Fiel	ds: 15 in, 0 filtered, 3 renamed, 15 out
Field	Filter	Field
debtinc -		debtinc 🖌
creddebt -		creddebt
othdebt -		othdebt
default -		default
\$B-default -		TAN
\$BP-default -		\$BP-default
\$B1-default -		Markov
\$BP1-default -		\$BP1-default
\$B2-default -		Markov-FS
\$BP2-default -		\$BP2-default
View current fields O View unus OK Cancel	ed field s	settings

Figure 238. Modifiez les noms des champs de modèle

Pour comparer l'exactitude des prédictions des modèles, vous pouvez créer un graphique de gain.

7. Liez un noeud de graphique d'évaluation au noeud Filtre et exécutez le noeud Graphique en utilisant les paramètres par défaut.

Le graphique montre que chaque type de modèle produit des résultats similaires ; mais, le modèle Markov est légèrement meilleur.



Figure 239. Evaluation de la précision du modèle

Pour vérifier la précision des prédictions de chaque modèle, vous pouvez utiliser le noeud Analyse à la place du graphique d'évaluation. Ceci affiche la précision en pourcentage à la fois pour les prédictions correctes et incorrectes.

8. Liez un noeud Analyse au noeud Filtre et exécutez le noeud Analyse en utilisant les paramètres par défaut.

Comme pour le graphique d'évaluation, ce noeud montre que le modèle Markov est légèrement meilleur dans ses prédictions ; cependant, le modèle Markov-FS n'est qu'à quelques points de pourcentage derrière ce modèle. Cela peut indiquer qu'il est peut-être préférable d'utiliser le modèle Markov-FS car celui-ci utilise un moins grand nombre d'entrées pour calculer ses résultats, ce qui diminue la durée de la collecte de données ainsi que la durée d'entrée et de traitement des données.

🔦 Analysis	of [default]				
<u> F</u> ile	🖹 <u>E</u> dit		8		@ ×
Analysis	Appotations				
/ andiyono	Annotations				
8 Collar	ose All 🤷 E	kpand All			
Result	s for output field	default			
i⊟⊡Ind	ividual Models				
	Comparing TAN	with defa	ult		
	Correct	565	80.71%		
	Wrong	135	19.29%		
	Total	700			
Þ	Comparing Mar	kov with d	efault		
	Correct	542	77.43%		
	Wrong	158	22.57%		
	Total	700			
Ē.	Comparing Mar	kov-FS wit	h default		
	Correct	542	77.43%		
	Wrong	158	22.57%		
	Total	700			
⊟- Agr	eement betweer	TAN Mai	kov Markov-F	S	
	Agree	603	86.14%		
	Disagree	97	13.86%		
	Total	700	th default		
	Comparing Agre	ement wi	00 75%		
	Wrong	00	16 25%		
	Total	603	10.2376		
	Total	003			
					OK

Figure 240. Analyse de la précision des modèles

Des explications sur les fondements mathématiques des méthodes de modélisation utilisées dans IBM SPSS Modeler sont présentées dans le *guide des algorithmes d'IBM SPSS Modeler*, disponible dans le répertoire *Documentation* du disque d'installation.

Sachez également que ces résultats sont établis uniquement d'après les données d'apprentissage. Pour évaluer comment le modèle peut se généraliser à d'autres données dans le monde réel, vous devez utiliser un noeud Partitionner destiné à contenir un sous-ensemble d'enregistrements à des fins de test et de validation.

Chapitre 18. Recyclage d'un modèle chaque mois (Réseau Bayésien)

Les réseaux Bayésiens permettent de créer un modèle de probabilité en combinant les preuves observées et enregistrées avec les connaissances réelles "de bon sens" pour établir la probabilité des occurences en utilisant des attributs apparemment sans lien.

Cet exemple utilise le flux *bayes_churn_retrain.str*, qui fait référence aux fichiers de données *telco_Jan.sav* et *telco_Feb.sav*. Ces fichiers sont accessibles dans le répertoire *Demos* de toute installation IBM SPSS Modeler. Vous pouvez y accéder à partir du groupe de programmes IBM SPSS Modeler du menu Démarrer de Windows. Le fichier *bayes_churn_retrain.str* se trouve dans le répertoire des *flux*.

Par exemple, supposons qu'un fournisseur de télécommunications souhaite connaître le nombre de clients qui partent à la concurrence (attrition). Si les données client historiques peuvent être utilisées pour prédire quels clients sont les plus susceptibles d'attrition, ces clients peuvent être ciblés afin de recevoir des remises ou des offres pour les dissuader de passer à un autre fournisseur de services.

Cet exemple examine l'utilisation de données d'attrition mensuelles existantes pour prédire quels clients sont susceptibles d'attrition et l'ajout des données du mois suivant pour affiner et recycler ce modèle.

Création du flux

1. Ajoutez un noeud source Statistics pointant vers telco_Jan.sav dans le dossier Demos.



Figure 241. Flux d'échantillons de réseau Bayésien

L'analyse précédente a démontré que plusieurs champs de données sont inutiles pour la prédiction de l'attrition. Ces champs peuvent être filtrés à partir de votre jeu de données afin de réduire la durée du traitement lors de la création et de l'évaluation de modèles.

2. Ajoutez un noeud Filtrer au noeud Source.

- **3**. Excluez tous les champs à l'exception des champs *address, age, churn, custcat, ed, employ, gender, marital, reside, retire,* et *tenure.*
- 4. Cliquez sur OK.

42 in, 31 filtered, 0 renamed, 11 o
Field
region
tenure
age
marital
address
income
ed
employ
retire
gender

Figure 242. Filtrage des champs inutiles

- 5. Ajoutez un noeud type au noeud Filtrer.
- 6. Ouvrez le noeud type et cliquez sur le bouton Lire les valeurs pour remplir la colonne Valeurs.
- 7. Pour que le noeud Evaluation puisse évaluer les valeurs True (vrai) et False (faux), définissez le niveau de mesure pour le champ *attrition* sur **Indicateur**, et son rôle sur **Cible**. Cliquez sur **OK**.

Type	Annotations				0
4- 00	🗪 🚺 🕨 Read Va	lues Clear	^r Values	Clear All Va	alues
Field	Measurement	Values	Missing	Check	Role
🖌 manitai	🍯 riag	170		NOTE	🗯 ii iput 🔽
🔉 address	🔗 Continuous	[0,55]		None	🔪 Input
ݤ ed	💑 Nominal	1,2,3,4,5		None	🔪 Input
> employ	🖉 Continuous	[0,47]		None	N Input
retire	🖁 Flag	1.0/0.0		None	> Input
gender	Nominal	0,1		None	> Input
> reside	🂑 Nominal	1,2,3,4,5,		None	N Input
> custcat	💑 Nominal	1,2,3,4		None	> Input
churn	🖁 Flag	1/0		None	O Target
View current OK Cance	t fields 🔘 View unu	sed field settin	gs		Apply Rese

Figure 243. Sélection du champ cible

Vous pouvez créer plusieurs types de réseaux Bayésiens ; mais, dans cet exemple, vous créez un modèle Tree Augmented Naïve Bayes (TAN). Celui-ci crée un réseau étendu qui vous permet d'inclure tous les liens possibles entre les variables de données, créant ainsi un modèle initial fiable.

- 8. Reliez un noeud de réseau Bayésien au noeud type.
- 9. Pour choisir le nom du modèle, dans l'onglet Modèle, sélectionnez **Personnalisé** puis saisissez Jan dans la zone de texte.
- 10. Pour la méthode d'apprentissage des paramètres, sélectionnez Ajustement de Bayes pour les petits calculs de cellules.
- 11. Cliquez sur **Exécuter**. Le nugget de modèle est ajouté au flux et également à la palette Modèles en haut à droite.

😨 Jan 🛛 🕅
Model name: O Auto O Custom Jan
Use partitioned data
☑ Build model for each split
To select fields manually, choose "Use custom settings" on the Fields tab
Partition:
Splits:
Continue training existing model
Structure type: 🔘 TAN 🔘 Markov Blanket
Include feature selection preprocessing step
Parameter learning method:
◎ Maximum likelihood ◎ Bayes adjustment for small cell counts
OK Run Cancel Apply Reset

Figure 244. Création d'un modèle Tree Augmented Naïve Bayes

- 12. Ajoutez un noeud source Statistics pointant vers telco_Feb.sav dans le dossier Demos.
- **13**. Attachez ce nouveau noeud source au noeud Filtrer (dans la boîte de dialogue d'avertissement, choisissez **Remplacer** pour remplacer la connexion au noeud source précédent).



Figure 245. Ajout des données du deuxième mois

- 14. Pour choisir le nom du modèle, dans l'onglet Modèle du noeud Réseau Bayésien, sélectionnez **Personnalisé** puis saisissez Jan-Feb dans la zone de texte.
- 15. Sélectionnez Poursuivre l'apprentissage du modèle existant.
- 16. Cliquez sur **Exécuter**. Le nugget de modèle remplace le nugget existant dans le flux mais il est également ajouté à la palette Modèles en haut à droite.

🙀 Jan-Feb		×
		0-0
Fields Model Expert	Analyze Annotations	
Model name:	🔘 Auto 💿 Custom	Jan-Feb
👿 Use partitioned data		
👿 Build model for each :	split	
To select fields manually	, choose "Use custom settin	gs" on the Fields tab
Partition:		-
Splits:		×
Continue training existi	ing model	
Structure type:	🔘 TAN 🔘 Markov Bla	nket
🔲 Include feature selecti	on preprocessing step	
Parameter learning metho	d:	
	🔘 Maximum likelihood 🄇	Bayes adjustment for small cell counts
OK 🕨 Run	Cancel	Apply <u>R</u> eset

Figure 246. Recyclage du modèle

Evaluation du modèle

Pour comparer les modèles, vous devez combiner les deux jeux de données.

1. Ajoutez un noeud Ajouter et liez-y les noeuds source telco_Jan.sav et telco_Feb.sav.

😡 Append		
Preview)		
Append 2 datase	ts	
Inputs Append Annotation	ns	
Match fields by: O Positi	on 🔘 Name 📃 Mato	h case
Preview of field matches ar	nd structure	
Output Field -	1[telco_Jan.sav:telco_Jan.s	2[telco_Feb.sav:telco_Feb
🔆 region	📿 region	🔆 region 🛛 🔺
🔆 tenure	🔆 tenure	🔆 tenure
🔷 age	🔷 age	🔷 age 🗾
🔆 marital	🔆 marital	🚫 marital
🔆 address	🔆 address	🚫 address
🛞 income	🛞 income	🛞 income
🔆 ed	🔷 ed	🔆 ed
🔆 employ	🚫 employ	🔆 employ 🔽
leskula fielde fram: 🙆 Main	data ant ank. 🔿 till datat-	
include neids from: SMain	uataset only 🥌 All datasets	
Tag records by including s	ource dataset in field Input	
OK Cancel		Apply Reset

Figure 247. Ajout des deux sources de données

- 2. Copiez les noeuds Filtrer et type précédents dans le flux et collez-les dans le canevas de flux.
- 3. Liez le noeud Ajouter au noeud Filtrer que vous venez de coller.



Figure 248. Collage des noeuds copiés dans le flux

Les nuggets des deux modèles de réseau Bayésien se trouvent dans la palette Modèles en haut à droite.

- 4. Double-cliquez sur le nugget de modèle Jan pour l'ajouter au flux et reliez-le au nouveau noeud copié type.
- 5. Liez le nugget de modèle Jan-Feb déjà dans le flux au nugget de modèle Jan.
- 6. Ouvrez le nugget de modèle Jan.



Figure 249. Ajout des nuggets au flux

L'onglet Modèle du nugget de modèle Réseau Bayésien est divisé en deux colonnes. La colonne de gauche contient un graphique de noeuds en réseau qui affiche la relation entre la cible et ses prédicteurs les plus importants, ainsi que la relation entre les prédicteurs.

La colonne de droite affiche soit l'*Importance des prédicteurs*, qui indique l'importance relative de chaque prédicteur pour l'estimation du modèle ou les *Probabilités conditionnelles*, qui contiennent la valeur de la probabilité conditionnelle de chaque valeur de noeud et pour chaque combinaison de valeurs dans ses noeuds parent.



Figure 250. Modèle Réseau Bayésien montrant l'importance des prédicteurs

Pour afficher les probabilités conditionnelles d'un noeud, cliquez sur ce noeud dans la colonne de gauche. La colonne de droite est mise à jour et affiche les détails appropriés.

Les probabilités conditionnelles sont affichées pour chaque intervalle de division des valeurs de données associé aux noeuds parent et aux noeuds frères.



Figure 251. Modèle Réseau Bayésien affichant les probabilités conditionnelles

- 7. Pour renommer les sorties de modèles et les rendre plus claires, liez un noeud Filtrer au nugget de modèle Jan-Feb.
- 8. Dans la colonne de droite *Champ*, remplacez le nom \$B-attrition par Jan et \$B1-attrition par Jan-Fév.

Filter		
	Field	ls: 15 in, 0 filtered, 2 renamed, 15 out
Field	Filter	Field
employ		employ 🖌
retire	\rightarrow	retire
gender	\rightarrow	gender
reside	\rightarrow	reside
custcat	\rightarrow	custcat
churn	\rightarrow	churn
\$B-churn	\rightarrow	Jan
\$BP-churn	\rightarrow	\$BP-churn
\$B1-churn	\rightarrow	Jan-Feb
\$BP1-churn	\rightarrow	\$BP1-churn
View current fields View unu OK Cancel	used field set	ttings

Figure 252. Modifiez les noms des champs de modèle

Pour vérifier l'exactitude des prédictions d'attrition de chaque modèle, utilisez un noeud Analyse qui permet d'afficher cette exactitude en termes de pourcentage, à la fois pour les prédictions correctes et incorrectes.

- 9. Reliez un noeud Analyse au noeud Filtrer.
- 10. Ouvrez le noeud Analyse, puis cliquez sur Exécuter.

Cela montre que les deux modèles ont des exactitudes semblables lors de la prédiction d'attrition.

🔦 Analysi	s of [churn]			
📦 File 🛛	🛓 Edit 🛛 🐻	<u>0</u> 14		0 ×
Analysis	Annotations			
8 Collaps	e All 🧛 Exp	and All		
Results - Indiv	for output field ch vidual Models Comparing Jan wi	urn th churn		
	Correct	771	77.1%	
	Wrong	229	22.9%	
	Total	1,000		
<u> </u>	Comparing Jan-Fe	b with chu	'n	
	Correct	765	76.5%	
	Wrong	235	23.5%	
	Total	1,000		
🖨 Agn	eement between	Jan Jan-Fel	0	
	Agree	882	88.2%	
	Disagree	118	11.8%	
	Total	1,000		
ė.	Comparing Agree	nent with c	hurn	
	Correct	710	80.5%	
	Wrong	172	19.5%	
	Total	882		
				ОК

Figure 253. Analyse de la précision des modèles

A la place du noeud Analyse, vous pouvez utiliser un graphique Evaluation pour comparer la précision des prédictions des modèles en créant un graphique de gain.

11. Ajoutez un noeud de graphique Evaluation au noeud Filtrer

et exécutez le noeud Graphique en utilisant ses paramètres par défaut.

Comme le noeud Analyse, ce graphique affiche que chaque type de modèle produit des résultats similaires ; mais, le modèle recyclé qui utilise les données des deux mois est légèrement meilleur car ses prédictions ont un plus haut niveau de fiabilité.



Figure 254. Evaluation de la précision du modèle

Des explications sur les fondements mathématiques des méthodes de modélisation utilisées dans IBM SPSS Modeler sont présentées dans le *guide des algorithmes d'IBM SPSS Modeler*, disponible dans le répertoire *Documentation* du disque d'installation.

Sachez également que ces résultats sont établis uniquement d'après les données d'apprentissage. Pour évaluer comment le modèle peut se généraliser à d'autres données dans le monde réel, vous devez utiliser un noeud Partitionner destiné à contenir un sous-ensemble d'enregistrements à des fins de test et de validation.

Chapitre 19. Campagne publicitaire (Réseau de neurones/Arbre C&RT)

Cet exemple s'appuie sur des données relatives à des gammes de produits destinés à la vente au détail et aux effets de la campagne publicitaire sur les ventes. (Ces données sont fictives.) Votre objectif, dans cet exemple, est de prévoir les effets des prochaines campagnes publicitaires. Comme dans l'exemple de surveillance d'état, le processus d'exploration de données se compose des étapes suivantes : exploration, préparation des données, apprentissage et tests.

Cet exemple utilise les flux nommés *goodsplot.str* et *goodslearn.str*, qui font référence aux fichiers de données nommés *GOODS1n* et *GOODS2n*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le flux *goodsplot.str* se situe dans le dossier des *flux*, tandis que le fichier *goodslearn.str* se situe dans le répertoire des *flux*.

Analyse des données

Chaque enregistrement comprend les éléments suivants :

- *Classe*. Type de produit.
- Coût. Prix unitaire.
- Campagne publicitaire. Somme consacrée à une campagne publicitaire particulière.
- Avant. Recettes avant la campagne publicitaire.
- Après. Recettes après la campagne publicitaire.

Le flux *goodsplot.str* contient un flux simple permettant d'afficher les données dans un tableau. Les deux champs relatifs aux recettes (*Avant* et *Après*) sont exprimés en termes absolus ; cependant, la valeur de l'augmentation des recettes après la campagne publicitaire (certainement due à celle-ci) vous sera probablement plus utile.

Promotions							
違 <u>F</u> ile	📄 Edit	🏷 Gene	erate 🧯		14 3	8 0 🗙	
Table A	Annotations						
	Class	Cost	Promotion	Before	After		
1	Confection	23.990	1467	114957	122762	-	
2	Drink	79.290	1745	123378	137097		
3	Luxury	81.990	1426	135246	141172		
4	Confection	74.180	1098	231389	244456		
5	Confection	90.090	1968	235648	261940		
6	Meat	69.850	1486	148885	156232		
7	Meat	100.1	1248	123760	128441		
8	Luxury	21.010	1364	251072	268134		
9	Luxury	87.320	1585	287043	310857		
10	Drink	26.580	1835	240805	272863		
11	Drink	65.230	1194	212406	227836		
12	Meat	79.820	1596	174022	181489		
13	Confection	41.390	1161	270631	283189		
14	Meat	36.820	1151	231281	235722		
15	Meat	44.050	1482	178138	185934		
16	Drink	84.620	1623	247885	278031		
17	Confection	51.820	1969	148597	165598		
18	Confection	90.080	1462	215102	228696		
19	Luxury	57.300	1842	246885	270082		
20	Drink	11.020	1370	164984	176802	-	

Figure 255. Effets de la campagne publicitaire sur les ventes du produit

goodsplot.str contient également un noeud pour calculer cette valeur exprimée sous la forme d'un pourcentage de la recette avant la campagne publicitaire, dans un champ appelé *Augmentation*, et affiche un tableau indiquant ce champ.

🖩 Promotions 📃 🗆 🔀							
じ <u>F</u> ile	📄 Edit	🏷 <u>G</u> ene	erate 🚺		14		0 ×
Table	Annotations						
	Class	Cost	Promotion	Before	After	Increase	
1	Confection	23.990	1467	114957	122762	6.789	4
2	Drink	79.290	1745	123378	137097	11.119	
3	Luxury	81.990	1426	135246	141172	4.382	
4	Confection	74.180	1098	231389	244456	5.647	
5	Confection	90.090	1968	235648	261940	11.157	
6	Meat	69.850	1486	148885	156232	4.935	
7	Meat	100.1	1248	123760	128441	3.782	
8	Luxury	21.010	1364	251072	268134	6.796	
9	Luxury	87.320	1585	287043	310857	8.296	
10	Drink	26.580	1835	240805	272863	13.313	
11	Drink	65.230	1194	212406	227836	7.264	
12	Meat	79.820	1596	174022	181489	4.291	
13	Confection	41.390	1161	270631	283189	4.640	Long and
14	Meat	36.820	1151	231281	235722	1.920	
15	Meat	44.050	1482	178138	185934	4.376	
16	Drink	84.620	1623	247885	278031	12.161	
17	Confection	51.820	1969	148597	165598	11.441	
18	Confection	90.080	1462	215102	228696	6.320	
19	Luxury	57.300	1842	246885	270082	9.396	
20	Drink	11.020	1370	164984	176802	7.163	-
							ОК

Figure 256. Augmentation des recettes après la campagne publicitaire

De plus, le flux affiche un histogramme de l'augmentation et un nuage de points de l'augmentation par rapport aux coûts de la campagne publicitaire, avec la catégorie de produits concernée.



Figure 257. Histogramme de l'augmentation des recettes

Le nuage de points fait apparaître que, pour chaque catégorie de produits, il existe une relation quasi-linéaire entre l'augmentation des recettes et les coûts engagés dans la campagne publicitaire. Il est donc probable qu'un arbre décisions ou un réseau de neurones puisse prévoir, avec une exactitude relativement fiable, l'augmentation des recettes à l'aide des autres champs disponibles.



Figure 258. Rapport augmentation des recettes/dépenses publicitaires

Apprentissage et tests

Le flux goodslearn.str forme un réseau de neurones et un arbre décision pour effectuer cette prévision d'augmentation des recettes.



Figure 259. Modélisation du flux goodslearn.str (prodappren)

Une fois que vous avez exécuté les noeuds de modèle et généré les véritables modèles, vous pouvez tester les résultats du processus d'apprentissage. Pour ce faire, connectez l'arbre de décisions et le réseau en série entre le noeud type et un nouveau noeud analyse, remplacez le fichier (de données) d'entrée par PR0D1n, puis exécutez le noeud analyse. A partir du résultat de ce noeud, notamment de la corrélation linéaire entre l'augmentation prévue et le résultat réel, vous remarquerez que les systèmes formés prévoient l'augmentation des recettes avec un niveau de fiabilité élevé.

Une étude plus poussée montrerait des cas où les systèmes formés commettent des erreurs assez importantes ; ces erreurs peuvent être identifiées par la représentation graphique de l'augmentation prévue par rapport à l'augmentation réelle. Dans ce cas, il vous suffirait de sélectionner les données déviantes à l'aide des graphiques interactifs de SPSS Modeler et, dans leurs propriétés, de rectifier la description des données et le processus d'apprentissage pour améliorer l'exactitude des prévisions.

Chapitre 20. Surveillance d'état (Réseau de neurones/C5.0)

Cet exemple se rapporte à la surveillance des informations de statut à partir d'un ordinateur, et à la difficulté à identifier et à prévoir les statuts de panne. Les données sont créées à partir d'une simulation fictive et se composent de plusieurs séries concaténées mesurées dans le temps. Chaque enregistrement rend compte du dernier état de la machine en indiquant les informations suivantes :

- *Heure*. Un entier.
- Puissance. Un entier.
- *Température*. Un entier.
- Pression. 0 si la situation est normale, 1 pour avertir d'un risque passager de pression.
- Temps de bon fonctionnement. Temps écoulé depuis la dernière intervention.
- Statut. Généralement 0. Il prend la valeur du code d'erreur (101, 202 ou 303) en cas d'erreur.
- *Résultat*. Le code d'erreur qui apparaît dans les séries temporelles ou 0 si aucune erreur ne se produit. (ces codes ne sont disponibles qu'a posteriori).

Cet exemple utilise les flux nommés *condplot.str* et *condlearn.str*, qui font référence aux fichiers de données *COND1n* et *COND2n*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Les fichiers *condplot.str* et *condlearn.str* sont disponibles dans le répertoire des *flux*.

A chaque série temporelle correspond une série d'enregistrements commençant par une période d'activité normale suivie par la période menant à la panne, comme l'indique le tableau suivant :

Temps	Puissance	Température	Pression	Temps de bon fonctionnement	Statut	Résultat
0	1059	259	0	404	0	0
1	1059	259	0	404	0	0
			•••			
51	1059	259	0	404	0	0
52	1059	259	0	404	0	0
53	1007	259	0	404	0	303
54	998	259	0	404	0	303
			•••			
89	839	259	0	404	0	303
90	834	259	0	404	303	303
0	965	251	0	209	0	0
1	965	251	0	209	0	0
51	965	251	0	209	0	0
52	965	251	0	209	0	0
53	938	251	0	209	0	101
54	936	251	0	209	0	101

				Temps de bon			
Temps	Puissance	Température	Pression	fonctionnement	Statut	Résultat	
208	644	251	0	209	0	101	
209	640	251	0	209	101	101	

Le processus suivant est commun à la plupart des projets d'exploration de données :

- Examinez les données permettant de déterminer les attributs utiles à la prévision ou à l'identification des états à connaître.
- Conservez ces attributs (s'ils existent) ou calculez-les, puis ajoutez-les aux données, si cela est nécessaire.
- Utilisez les données obtenues pour l'apprentissage des règles et des réseaux de neurones.
- Testez les systèmes formés à l'aide de données de test indépendantes.

Analyse des données

Le fichier *condplot.str* illustre la première partie du processus. Il contient un flux qui reproduit plusieurs graphiques. Si les séries temporelles de température et de puissance contiennent des motifs visibles, vous pouvez distinguer les différentes conditions d'erreurs imminentes, voire prévoir le moment où l'erreur se produira. Pour la température et la puissance, le flux ci-dessous trace les séries temporelles associées aux trois codes d'erreur différents sur des graphiques individuels, produisant six graphiques. Les noeuds Sélectionner séparent les données associées aux différents codes d'erreur.





Les résultats de ce flux sont indiqués dans cette figure.



Figure 261. Température et puissance dans le temps

Les graphiques affichent clairement les motifs pouvant différencier jusqu'à 202 erreurs (de l'erreur 101 à l'erreur 303). Ces 202 erreurs rendent compte de l'augmentation de la température et de la fluctuation de la puissance dans le temps, ce qui n'est pas le cas des autres erreurs. Cependant, les motifs différenciant les erreurs 101 des erreurs 303 ne sont pas aussi explicites. Ces deux erreurs montrent une température régulière et une baisse de puissance, mais cette baisse est plus accentuée quand l'erreur 303 apparaît.

Ces graphiques permettent de constater que le changement et le degré de fluctuation de la température et de la puissance sont déterminants pour la prévision et l'identification des pannes. Ces attributs doivent donc être ajoutés aux données avant l'application des systèmes d'apprentissage.

Préparation des données

A partir des résultats de l'exploration de données, le flux *condlearn.st* calcule les données appropriées et apprend à prévoir les pannes.



Figure 262. Flux condlearn (étatappren)

Le flux utilise plusieurs noeuds dériver pour préparer les données en vue de la modélisation.

- Noeud de type Délimité. Lit le fichier de données ETAT1n.
- Noeud dériver Avertissements de pression. Calcule le nombre d'avertissements de pression passagers. Réinitialisez cette valeur lorsque l'heure revient sur la valeur 0.
- Noeud dériver AugTemp. Calcule le taux de changement passager de température à l'aide de @DIFF1.
- Noeud dériver AugPuiss. Calcule le taux de changement passager de la puissance à l'aide de @DIFF1.
- Noeud dériver FluxPuiss. Il s'agit d'un indicateur avec la valeur True (vraie) si les variations de la puissance sont totalement différentes entre le dernier enregistrement et celui-ci, c'est-à-dire entre un pic de puissance et une baisse importante.
- Noeud dériver EtatPuiss. L'état est *Stable* au début, puis devient *Fluctuant* lorsque deux flux de puissance successifs sont détectés. Il revient sur *Stable* uniquement lorsqu'il n'y a pas eu de flux de puissance pendant cinq intervalles de temps ou lorsque *Temps* est réinitialisé.
- ModifPuiss. Moyenne de AugPuiss sur les cinq derniers intervalles de temps.
- ModifTemp. Moyenne de AugTemp sur les cinq derniers intervalles de temps.
- Noeud Abandonner (sélectionner). Supprime le premier enregistrement de chaque série temporelle pour éviter des écarts importants (inappropriés) de *Puissance* et de *Température* entre chaque enregistrement.
- Noeud Abandonner les champs. Réduit les champs des enregistrements à *Temps de bon fonctionnement*, *Statut*, *Résultat*, *Avertissements de pression*, *EtatPuiss*, *ModifPuiss* et *ModifTemp*.
- **Type**. Définit le rôle du *Résultat* en tant que **Cible** (le champ à prédire). En outre, définit le niveau de mesure de *Résultat* en tant que **Nominal**, *Avertissements de pression* en tant que **Continu** et *EtatPuiss* en tant que **Indicateur**.

Apprentissage

L'exécution du flux dans *condlearn.str* permet l'apprentissage de la règle C5.0 et du réseau de neurones. Le temps d'apprentissage du réseau peut être long, mais vous pouvez interrompre le processus avant la fin pour enregistrer un réseau produisant des résultats satisfaisants. Une fois l'apprentissage terminé, l'onglet Modèles en haut à droite dans la fenêtre des gestionnaires clignote pour vous avertir que deux nuggets ont été créés : l'un représente le réseau de neurones et l'autre, la règle.



Figure 263. Gestionnaire de modèles avec des nuggets de modèle

Les nuggets de modèle sont aussi ajoutés au flux existant, ce qui nous permet de tester le système ou d'exporter les résultats du modèle. Dans cet exemple, nous allons tester les résultats du modèle.

Test

Les nuggets de modèle sont ajoutés au flux, tous deux étant connectés au noeud type.

- 1. Repositionnez les nuggets comme indiqué, de sorte que le noeud type se connecte au nugget Réseau de neurones, qui se connecte au nugget C5.0.
- 2. Reliez un noeud Analyse au nugget C5.0.
- **3**. Editez le noeud source d'origine pour lire le fichier *ETAT2n* (au lieu de *ETAT1n*), car *ETAT2n* contient des données de test non affichées.



Figure 264. Test du réseau formé

4. Ouvrez le noeud Analyse, puis cliquez sur Exécuter.

Ceci génère des résultats reflétant l'exactitude du réseau formé et de la règle.

Chapitre 21. Classification des clients de services de télécommunications (analyse discriminante)

L'analyse discriminante est une technique statistique de classification des enregistrements sur la base des valeurs des champs d'entrée. Excepté le fait qu'elle utilise un champ cible catégoriel et non pas numérique, cette régression est semblable à la régression linéaire.

Par exemple, supposons qu'un fournisseur de télécommunications ait segmenté sa base de clientèle par modèles d'utilisation de service, classant ses clients en quatre groupes. Si les données démographiques peuvent être utilisées pour prévoir les affectations de groupes, vous pouvez personnaliser les offres pour les clients éventuels.

Cet exemple utilise le flux *telco_custcat_discriminant.str*, qui fait référence au fichier de données *telco.sav*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *telco_custcat_discriminant.str* se trouve dans le répertoire des *flux*.

Cet exemple est axé sur l'utilisation des données démographiques dans le but de prévoir des modèles d'utilisation. Le champ cible *custcat* possède quatre valeurs possibles qui correspondent aux quatre groupes de clients suivants :

Valeur	Libellé
1	Basic Service
2	E-Service
3	Plus Service
4	Total Service

Création du flux

1. Définissez en premier lieu des propriétés de flux pour afficher les libellés de variable et de valeur dans la sortie. Dans les menus, sélectionnez :

Fichier > Propriétés du flux... > Options > Général

2. Sélectionnez Afficher les libellés de champ et de valeur dans le résultat et cliquez sur OK.

😵 telco_custcat_discriminant 🛛 🛛 🕅									
									0
Options Messages	Parameters	Deployment	Script	Globals	Search	Comments	Annotations		
Select a setting:									
General	These are g for all your s	eneral setting: streams.	s that app	ly to the d	current str	eam. Click Sa	ave As Default to u	use these settings	as the default
Date/Time									
Number formats	Decimal syn	nbol:		P	eriod (.)	-			
Optimization	<u>G</u> rouping sy	/mbol:		N	lone	-			
Logging and Status	Encoding:			s	ystem def	ault 🔻			
Layout	Ruleset Eva	luation:		V	'oting 💌	1			
	— Maximum п.	umber of rows	to show	in Data P	review:		10 4		
				Galala			250		
		n members to	rnominali	ieids		Ļ	250		
	lumit set	size for Koho	onen and l	<-Means	modeling	L	20 🔫		
	C Re <u>f</u> resh	source node:	s on exect	ution					
	📝 Display	field and value	e labels in	output					
									Save As Default
OK Cancel									Apply Reset

Figure 265. Propriétés du flux

3. Ajoutez un noeud source de fichier Statistics pointant vers telco.sav dans le dossier Demos.



Figure 266. Flux d'échantillons permettant de classifier les clients par analyse discriminante

a. Ajoutez un noeud type et cliquez sur **Lire les valeurs**, en vous assurant que tous les niveaux de mesure sont correctement paramétrés. Par exemple, la majorité des champs avec des valeurs 0 et 1 peuvent être considérés comme des champs indicateurs.

Type								
~	Pormat Annotations Pread Values Clear Values Clear All Values							
Field -	Measurem	ent	Values	Missing		Check	R	ole
🔆 gender	💑 Nominal		0,1		Nor	ne	> Inp	ut 📥
🚫 reside	🖉 Continuous	()	[1,8]		Nor	ne	🔪 Inp	ut
🔿 tollfree	🎖 Flag		1/0		Nor	ne	🔪 Inp	ut
📿 equip	🎖 Flag		1/0		Nor	ne	🔪 Inp	ut
📿 callcard	Flag		1./0		N	<default:< td=""><td>></td><td></td></default:<>	>	
父 wireless	Flag		1/0		NZ			
🛞 longmon	🖉 Continuou	Se	ect All			Continuo	us	
😤 tollmon 🖉 Continuou s			Select None			Categorical		
View current	Select Fields			Flag				
OK	Co	py ste Special	Ctrl+C . Ctrl+V		Nominal Ordinal	h	Reset	

Figure 267. Configuration du niveau de mesure pour plusieurs champs

Conseil : Pour modifier les propriétés de plusieurs champs contenant des valeurs similaires (telles que 0/1), cliquez sur l'en-tête de colonne *Valeurs* afin de trier les champs en fonction de cette valeur. Maintenez la touche Maj enfoncée tout en utilisant la souris ou les touches fléchées pour sélectionner tous les champs à modifier. Cliquez ensuite sur la sélection avec le bouton droit de la souris pour modifier le niveau de mesure ou les autres attributs des champs sélectionnés.

Veuillez noter que puisqu'il est plus correct de considérer le *sexe* comme un champ avec un ensemble de deux valeurs plutôt que comme un indicateur, laissez sa valeur de mesure sur **Nominal**.

b. Définissez le rôle du champ *custcat* sur **Cible**. Le rôle de tous les autres champs doit être défini sur **Entrée**.

Type	Preview) Annotations	duras Class	Values				
Field -	Measurement	Values	Missing	Check	Role		
epili	nag	170	3	NUTE	a input		
P logiong	Continuous	[-0.10536		None	Input		
P logtoll	Continuous	[1.74919		None	> Input		
ngequi 🖉	Continuous	[2.73436		None	> Input		
ngcard 🦉	Continuous	[1.01160		None	🔪 Input		
🤣 logwire	🔗 Continuous	[2.70136		None	🔪 Input		
🤔 Ininc	🔗 Continuous	[2.19722		None	🔪 Input		
🔆 custcat	💑 Nominal	1,2,3,4		None	🔘 Target		
🔆 churn	💑 Nominal	0,1		None	🔪 Input		
Immo Continuous [2.19722 None Input Custcat Nominal 1,2,3,4 None Input Churn Nominal 0,1 None Input Imput View current fields View unused field settings							

Figure 268. Définition du rôle de champ

Cet exemple étant axé sur les données démographiques, utilisez un noeud Filtrer pour n'inclure que les champs pertinents (*region, age, marital, address, income, ed, employ, retire, gender, reside* et *custcat*). Les autres champs peuvent être exclus pour cette analyse.

Demographic		
7.	Fields:	42 in, 31 filtered, 0 renamed, 11 o
Field -	Filter	Field
region	\rightarrow	region
tenure	× >	tenure
age	\rightarrow	age
marital	\rightarrow	marital
address	\rightarrow	address
income	\rightarrow	income
ed	\rightarrow	ed
employ	\rightarrow	employ
retire	\rightarrow	retire
gender	\rightarrow	gender
View current fields View OK Cancel	unused field	settings

Figure 269. Filtrage des champs démographiques

(Vous pouvez également paramétrer le rôle sur **Aucun** pour ces champs plutôt que de les exclure, ou sélectionner les champs que vous souhaitez utiliser dans le noeud de modélisation.)

4. Dans le noeud discriminant, cliquez sur l'onglet Modèle et sélectionnez la méthode Pas à pas.

😡 custcat			
			0
Fields Model	Expert Analyze Anr	notations	
Model name:	🔘 Auto 🔇	🕽 Custom	
👿 Use partitione	ed data		
👿 Build model fo	or each split		
Method: Stepwis	ie T		
ок 🕨 і	Run Cancel		Apply Reset

Figure 270. Choix des options de modèle

5. Dans l'onglet Expert, paramétrez le mode sur Expert et cliquez sur Sortie.

6. Sélectionnez **Tableau récapitulatif**, **Carte territoriale** et **Récapitulatif des étapes** dans la boîte de dialogue Sorties avancées puis cliquez sur **OK**.

😡 Discriminant: Advance	ed Output 🛛 🔀					
Statistics						
Descriptives:	Matrices:					
Means	Within-groups correlation					
🔲 Univariate ANOVAS	📕 Within-group covariance					
🔲 Box's M	📕 Separate-groups covariance					
Function Coefficients:	Total covariance					
E Fisher's						
🔄 Unstandardized						
Classification						
🔲 Casewise results	Plots:					
Limit cases to first:	10 🤤 🗹 Territorial map					
👿 Summary table	Combined-groups					
E Leave-one-out classification	on 📃 Separate-groups					
Stepwise						
👿 Summary of Steps						
F for pairwise distances						
Cancel Help						

Figure 271. Choix des options de sortie

Examen du modèle

1. Cliquez sur **Exécuter** pour créer le modèle qui est ajouté au flux et à la palette Modèles en haut à droite. Pour afficher ses détails, double-cliquez sur le nugget de modèle du flux.

L'onglet Récapitulatif affiche (entre autres) la cible et la liste complète des entrées (champs prédicteurs) à examiner.

😡 custca	t				
	iie <u>F</u> ile	🏷 Gen	erate 🚺	Preview	
Model Ad	dvanced	Settings se All	Summary	Annotations	
Analysi	is get Second get get get get get get get get	cat al ess me er ler le			

Figure 272. Récapitulatif du modèle avec champs cible et champs d'entrée

Pour des détails sur les résultats de l'analyse discriminante :

- 2. Cliquez sur l'onglet Avancé.
- **3**. Cliquez sur le bouton « Lancer dans un navigateur externe » (juste en-dessous de l'onglet Modèle) pour afficher les résultats dans votre navigateur Web.

Etude des résultats de l'utilisation de l'analyse discriminante pour classifier les clients de services de télécommunications

Analyse discriminante étape par étape

Variables Not in the Analysis

Step		Tolerance	Min. Tolerance	F to Enter	Wilks' Lambda
0	Age in years	1.000	1.000	7.521	.978
	Marital status	1.000	1.000	3.500	.990
	Years at current address	1.000	1.000	8.433	.975
	Household income in thousands	1.000	1.000	6.689	.980
	Level of education	1.000	1.000	61.454	.844
	Years with current employer	1.000	1.000	16.976	.951
	Retired	1.000	1.000	3.005	.991
	Gender	1.000	1.000	.373	.999
	Number of people in household	1.000	1.000	3.976	.988
1	Age in years	.980	.980	6.125	.829
	Marital status	.999	.999	3.803	.834
	Years at current address	.983	.983	8.487	.823
	Household income in thousands	.989	.989	6.022	.829
	Years with current employer	.953	.953	14.933	.807
	Retired	.992	.992	1.432	.840
	Gender	1.000	1.000	.358	.843
	Number of people in household	1.000	1.000	3.967	.834
2	Age in years	.563	.548	.352	.807
	Marital status	.999	.952	3.903	.798
	Years at current address	.798	.773	2.913	.800
	Household income in thousands	.689	.664	.634	.806
	Retired	.927	.891	.528	.806
	Gender	.998	.951	.391	.807
	Number of people in household	.979	.934	4.841	.796
3	Age in years	.535	.535	.252	.795
	Marital status	.605	.593	1.507	.792
	Years at current address	.776	.771	3.514	.787
	Household income in thousands	.688	.657	.687	.794
	Retired	.917	.880	.353	.795
	Gender	.997	.931	.395	.795

Figure 273. Variables absentes de l'analyse

Lorsque vous disposez de nombreux prédicteurs, la méthode détaillée étape par étape peut être utile car elle sélectionne automatiquement les "meilleures" variables à utiliser dans le modèle. La méthode détaillée étape par étape commence par un modèle qui n'inclut aucun des prédicteurs. A chaque étape, le prédicteur doté de la valeur *F-to-enter* la plus importante, supérieure aux critères d'entrée (par défaut, 3,84), est ajoutée au modèle.

Les variables qui ne sont toujours pas prises en compte dans l'analyse lors de la dernière étape ont toutes des valeurs *F-to-enter* inférieures à 3,84. Aucune variable supplémentaire n'est donc ajoutée.

Variables in the Analysis

Step		Tolerance	F to Remove	Wilks' Lambda
1	Level of education	1.000	61.454	
2	Level of education	.953	59.108	.951
	Years with current employer	.953	14.933	.844
3	Level of education	.951	60.046	.940
	Years with current employer	.934	15.824	.834
	Number of people in household	.979	4.841	.807

Figure 274. Variables comprises dans l'analyse

Cette table affiche les statistiques des variables qui figurent dans l'analyse à chaque étape. La *tolérance* est la proportion de la variance d'une variable non justifiée par les autres variables indépendantes de l'équation. Une variable ayant une très faible tolérance n'apporte que peu d'informations à un modèle et peut générer des problèmes de calcul.

Les valeurs *F-to-remove* sont utiles pour décrire ce qui se passe si une variable est supprimée du modèle actuel (les autres variables étant conservées). La valeur *F-to-remove* de la variable entrante est identique à la valeur *F-to-enter* de l'étape précédente (affichée dans la table Variables absentes de l'analyse).

Avertissement concernant les méthodes détaillées étape par étape

Les méthodes détaillées étape par étape sont pratiques mais connaissent leurs limites. Sachez que, étant donné que les méthodes détaillées étape par étape sélectionnent des modèles uniquement sur la base du mérite statistique, elles risquent de choisir des prédicteurs qui n'ont aucune *signification pratique*. Si vous connaissez bien les données et que vous avez des attentes particulières en ce qui concerne les prédicteurs importants, utilisez ces connaissances et évitez les méthodes détaillées étape par étape. Si, à l'inverse, vous avez de nombreux prédicteurs et que vous ne savez pas par où commencer, une analyse étape par étape et un ajustement du modèle sélectionné sont préférables à une absence complète de modèle.

Vérification de la qualité de l'ajustement

Eigenvalues								
Function	Eigenvalue	% of Variance	Cumulative %	Canonical Correlation				
1	.198ª	80.2	80.2	.407				
2	.048ª	19.4	99.6	.214				
3	.001ª	.4	100.0	.031				

a. First 3 canonical discriminant functions were used in the analysis.

Figure 275. Valeurs propres

Presque toute la variance expliquée par le modèle est liée aux deux premières fonctions discriminantes. Trois fonctions sont automatiquement ajustées. Cependant, du fait de sa valeur propre minime, vous pouvez ignorer la troisième en toute sécurité.

Wilks' Lambda

Test of Function(s)	Wilks' Lambda	Chi-square	df	Sig.
1 through 3	.796	227.345	9	.000
2 through 3	.953	47.486	4	.000
3	.999	.929	1	.335

Figure 276. lambda de Wilk

Le lambda de Wilk confirme que seules les deux premières fonctions sont utiles. Pour chaque ensemble de fonctions, cela permet de tester l'hypothèse selon laquelle les moyennes des fonctions répertoriées sont égales dans tous les groupes. Le test de la fonction 3 a une valeur de signification supérieure à 0,10. Par conséquent, cette fonction contribue peu au modèle.

Matrice de structure

Structure Matrix

	Function			
	1	2	3	
Level of education	.966*	090	244	
Years with current employer	182	.964*	193	
Age in years ^b	162	.598*	285	
Household income in thousands ^b	.109	.514*	190	
Years at current address ^b	151	.394*	214	
Retired ^b	108	.230*	137	
Gender ^b	.008	.054*	.009	
Number of people in household	.232	.097	.968*	
Marital status ^b	.132	.134	.600*	

Pooled within-groups correlations between discriminating variables and standardized canonical discriminant functions Variables ordered by absolute size of correlation within function.

*. Largest absolute correlation between each variable and any discriminant function

b. This variable not used in the analysis.

Figure 277. Matrice de structure

Lorsque plusieurs fonctions discriminantes existent, un astérisque (*) marque la corrélation absolue la plus importante de chaque variable avec l'une des fonctions canoniques. Dans chaque fonction, ces variables marquées sont ensuite triées en fonction de l'importance de la corrélation.

- La variable *Level of education* est la plus fortement corrélée avec la première fonction et est la seule variable la plus fortement corrélée avec cette fonction.
- Les variables *Years with current employer, Age in years, Household income in thousands, Years at current address, Retired* et *Gender* sont les plus fortement corrélées avec la deuxième fonction, bien que *Gender* et *Retired* soient plus faiblement corrélées que les autres. Les autres variables marquent cette fonction en tant que fonction de "stabilité".
- Les variables *Number of people in household* et *Marital status* sont les plus fortement corrélées avec la troisième fonction discriminante. Cependant, cette fonction discriminante étant sans intérêt, ces prédicteurs sont quasiment inutiles.

Carte territoriale

			Territor	rial Map					
_	(As	ssuming	g all func	tions bu	it the firs	t two an	e zero)		
Canonical L	Discrimina	ant							
-4.0	-3.0	-20	-1.0	0	1.0	20	3.0	4.0	
-4.0		-2.0	-1.0	.0	+	2.0	0.0	4.0	
4.0 +	1	1	1	3	4	1	+	1	
1				34			1		
1				34			1		
T				34			I		
T				34			I		
L				34			1		
3.0 +	+	+	+	+	34 +	+	+	+	
I				34			I		
T				34			1		
T				34			1		
I.				34			1		
I				34			1		
2.0 +	+	+	+	+ 3	4 +	+	+	+	
L				34			1		
1				34			I		
T				34			I.		
1			0	34			1		
T			(34			1		
1.0 +	+	+	+	+ 34	+	+	+	+	
I			З	324			1		
1			32	224			1		
L			32	2 24			1		
1			* 32	2 24			L		
1			32	24			1		
.0 +	+	+	+ 30	33332	*24 * +	+	+	+	
1		00000	001111	11112	24			1	
	000	00000	11111		10 04			1	
1	000	000011	11		10.04			1	
10000	2222111	1111			10.04			I I	
-10+11	11111	цці т	L 1		12 24	т	<u>т</u>	ц. Ц.	-
1.0 111		1	1	124	1224			1	1
i.				14			í.		
i.				14			Î		
i.				14			Î.		
i				14			i.		
-2.0 +	+	+	+	+	14	+	+	+	
1				14			L		
L.				14			Ī.		
L				14			L		
L				14			1		
L				14	ł		1		
-3.0 +	+	+	+	+	+ 14	+	+	+	
1				1	4		I		
L				1	4		1		
L				-	4		L		
L					14		1		
L					14		L		
-4.0 +					14		+		
+	+	+	+-		+	+	+	+	+
-4.0	-3.0	-2.0	-1.0	.0	1.0	2.0	3.0	4.0	
Symb	iols used	in territ	orial map)					
Symb	iol Group	b Labe							
1	1 Ba	sic serv	ice						
2	2 E-s	service							
3	3 Plu	s servic	e						
4	4 Tot	al servio	ce						
*	India	cates a	group ce	entroid					

Figure 278. Carte territoriale

La carte territoriale peut vous aider à étudier les relations entre les groupes et les fonctions discriminantes. Associée aux résultats de la matrice de structure, elle donne une interprétation graphique des relations entre prédicteurs et groupes. La première fonction, représentée sur l'axe horizontal, sépare le

groupe 4 (clients *Total service*) des autres. Etant donné que la variable *Level of education* est fortement corrélée avec la première fonction de manière positive, les clients *Total service* sont, en règle générale, ceux qui ont le niveau d'éducation le plus élevé. La deuxième fonction sépare les groupes 1 et 3 (clients *Basic service* et *Plus service*). En règle générale, les clients *Plus service* ont travaillé plus longtemps et sont plus âgés que les clients *Basic service*. Les clients *E-service* ne se distinguent pas nettement des autres, bien que la carte laisse penser qu'ils ont tendance à avoir un niveau d'éducation important et une expérience professionnelle moyenne.

En général, l'exactitude des centroïdes de groupe, marqués par des astérisques (*), par rapport aux lignes territoriales suggère que la séparation entre tous les groupes n'est pas très importante.

Seules les deux premières fonctions discriminantes sont tracées. Cependant, étant donné que la troisième fonction est relativement non significative, la carte territoriale fournit une vue complète du modèle discriminant.

Résultats de la classification supervisée

		Pre					
	Customer category	Basic service	E-service	Plus service	Total service	Total	
Original	Count	Basic service	125	11	61	69	266
		E-service	49	15	58	95	217
		Plus service	102	14	112	53	281
		Total service	40	16	37	143	236
	%	Basic service	47.0	4.1	22.9	25.9	100.0
		E-service	22.6	6.9	26.7	43.8	100.0
		Plus service	36.3	5.0	39.9	18.9	100.0
		Total service	16.9	6.8	15.7	60.6	100.0

Classification Results^a

a. 39.5% of original grouped cases correctly classified.

Figure 279. Résultats de la classification supervisée

Le lambda de Wilk vous permet de savoir que votre modèle permet d'obtenir des résultats pertinents, mais vous devez étudier les résultats de classification supervisée afin de déterminer à quel point ces résultats sont pertinents. D'après les données observées, le modèle "nul" (c'est-à-dire, un modèle sans prédicteurs) classifie tous les clients dans le groupe modal, *Plus service*. Par conséquent, le modèle nul serait correct 281/1 000 = 28,1 % du temps. Votre modèle obtient 11,4 % de plus, soit 39,5 % des clients. Votre modèle identifie particulièrement bien les clients *Total service*. Toutefois, il fonctionne très mal pour la classification des clients *E-service*. Vous devrez chercher un autre prédicteur pour séparer ces clients.

Récapitulatif

Vous avez créé un modèle discriminant qui classe les clients dans l'un des quatre groupes d'" utilisation de service " prédéfinis, en fonction des informations démographiques collectées auprès de chacun des clients. Grâce à la matrice de structure et à la carte territoriale, vous avez identifié les variables les plus utiles à la segmentation de votre clientèle. Enfin, les résultats de classification supervisée montrent que le modèle n'est pas très performant en ce qui concerne la classification des clients *E-service*. Davantage de recherches sont nécessaires pour déterminer une autre variable de prédicteur classant mieux ces clients.

Néanmoins, en fonction de ce que vous cherchez à prévoir, le modèle peut parfaitement correspondre à vos besoins. Par exemple, si l'identification des clients *E-service* ne vous intéresse pas, le modèle peut s'avérer assez précis pour vous. Cela peut être le cas lorsque le service en ligne est un produit d'appel qui n'engrange que peu de bénéfices. Si, par exemple, votre retour sur investissement le plus élevé provient des clients *Plus service* ou *Total service*, il est possible que le modèle vous fournisse les informations nécessaires.

Sachez également que ces résultats sont établis uniquement d'après les données d'apprentissage. Pour évaluer comment le modèle peut se généraliser à d'autres données, vous pouvez utiliser un noeud Partitionner destiné à contenir un sous-ensemble d'enregistrements à des fins de test et de validation.

Vous trouverez des explications sur le fondement mathématique des méthodes de modélisation utilisées dans IBM SPSS Modeler dans le Guide des algorithmes IBM SPSS Modeler. Celui-ci est disponible dans le répertoire *Documentation* du disque d'installation.
Chapitre 22. Analyse de données de survie avec censure par intervalle (modèles linéaires généralisés)

Lors de l'analyse de données de survie avec censure par intervalle (c'est-à-dire lorsque l'heure exacte de l'événement d'intérêt n'est pas connue, la seule donnée connue étant qu'il a eu lieu au cours d'un intervalle donné), l'application du modèle de Cox aux risques des événements sur des intervalles aboutit à un modèle de régression log-log complémentaire.

Des informations partielles, issues d'une étude visant à comparer l'efficacité de deux thérapies dans la prévention de la réapparition des ulcères, sont rassemblées dans le fichier *ulcer_recurrence.sav*. Ce jeu de données a été présenté et analysé ailleurs¹. A l'aide des modèles linéaires généralisés, vous pouvez répliquer les résultats pour les modèles de régression log-log complémentaires.

Cet exemple utilise le flux nommé *ulcer_genlin.str*, qui fait référence au fichier de données *ulcer_recurrence.sav*. Le fichier de données se trouve dans le dossier *Demos* et le fichier de flux dans le sous-dossier *streams*.

Création du flux

1. Ajoutez un noeud Statistics pointant vers *ulcer_recurrence.sav* dans le dossier *Demos*.



Figure 280. Flux d'échantillons relatif à la prévision de la réapparition des ulcères

2. Dans l'onglet Filtre du noeud source, excluez *id* et *time*.

^{1.} Collett, D. 2003. Modelling survival data in medical research, 2 ed. Boca Raton : Chapman & Hall/CRC.

🕜 ulcer	_recurrence.sav		×
	Preview Refr \$CLEO_DEMOSAlcer_rec	esh urrence.sav	0-0
Data Fi	tter Types Annotations		
7.	→ ₩		Fields: 6 in, 2 filtered, 0 renamed, 4 out
	Field -	Filter	Field
id		→	id
age		\rightarrow	age
duration		\rightarrow	duration
treatment		\rightarrow	treatment
time		→	time
result		\rightarrow	result
© View	current fields 🔘 View (unused field settings	Apply

Figure 281. Filtrage des champs superflus

- **3**. Dans l'onglet Types du noeud source, définissez le rôle du champ *résultats* sur **Cible** et son niveau de mesure sur **Indicateur**. Un résultat de 1 indique que l'ulcère est réapparu. Le rôle de tous les autres champs doit être défini sur **Entrée**.
- 4. Cliquez sur Lire les valeurs pour instancier les données.

\ - 000	🗪 🚺 🕨 Read Values	s Clear	Values 🚶	Clear All Valu	ies
Field 🗁	Measurement	Values	Missing	Check	Role
> age	🔗 Continuous	[23,76]		None	🔪 Input
duration	- Ordinal	1,2		None	🔪 Input
> treatment	al Nominal	0,1		None	🔪 Input
> result	🎖 Flag	1/0		None	🔘 Target

Figure 282. Définition du rôle de champ

5. Ajoutez un noeud Re-trier et spécifiez *duration, treatment* et *age* comme ordre des entrées. Il s'agit de l'ordre selon lequel les champs sont entrés dans le modèle ; il vous aide à répliquer les résultats de Collett.

🛛 Fiel	d Reorder		
	Preview		
Reorde	Annotations		
🔘 Custi	om Order	O Automatic Sort	
Туре:	Name: 🔺 💌 Storage	e: 🔺 🔻	
Туре	Field	Storage	
	······[other fields] ·····		
	duration	🚫 Integer	
8	treatment	🚫 Integer	
	age	🚫 Integer	^
			Ŧ
Clear I Note: Fi	Inused	not reordered.	
ок	Cancel		Apply Reset

Figure 283. Réorganisation des champs afin qu'ils soient entrés dans le modèle de la façon souhaitée

- 6. Reliez un noeud Modèles linéaires généralisés au noeud source. Dans le noeud Modèles linéaires généralisés, cliquez sur l'onglet **Modèle**.
- 7. Sélectionnez Premiers (valeur la plus faible) comme catégorie de référence pour la cible. Ce choix indique que la seconde catégorie est l'événement d'intérêt et que son effet sur le modèle se situe dans l'interprétation des estimations des paramètres. Un prédicteur continu avec un coefficient positif indique une probabilité de réapparition accrue, avec des valeurs croissantes du prédicteur ; les catégories d'un prédicteur nominal avec d'importants coefficients indiquent une probabilité de réapparition accrue par rapport aux autres catégories de l'ensemble.

🜍 result	
	0
Fields Model Expert Analyze Annotations	
Model name: O Auto Custom	
☑ Use partitioned data	
Build model for each split	
Model type: Main effects only Main effects and all two-way interactions	
Offset:	
Variable	
Offset field:	
O Fixed value	
Value: 0.0 🗢	
Base category for flag target: First (Lowest) 🔽	
☑ Include intercept in model	
OK Run Cancel	Apply Reset

Figure 284. Choix des options de modèle

- 8. Cliquez sur l'onglet Expert et sélectionnez Expert pour activer les options de modélisation expert.
- 9. Sélectionnez **Binomial** pour la distribution et **Log-log complémentaire** pour la fonction de lien.
- 10. Sélectionnez **Valeur fixe** comme méthode d'estimation du paramètre d'échelle et conservez la valeur par défaut 1,0.
- 11. Sélectionnez l'ordre des catégories des facteurs **Décroissant**. Cela signifie que la première catégorie de chaque facteur constitue la catégorie de référence. L'effet de cette sélection sur le modèle porte sur l'interprétation des estimations de paramètre.

😡 result		×
		0
Fields Model Expert Analy	ze Annotations	
Mode: 🔘 Simple 🍥 Expert		
Target Field Distribution and Link	Function	
The distribution that you choose	determines which	link functions are available.
Distribution: Binomial		Parameters
		Parameter for negative binomial:
		Specify value Value: 1.0
		Stimate
		Parameter for Tweedle: 1.5 🚔
Link function: Complementary lo	ig-log	Power: 0.0
Method and iteration settings are n	ot available if Distr	ribution = Normal and Link
Function = Identity.		
Parameter Estimation		
Method: Hybri	4 🔻	Maximum Fisher scoring iterations: 1
Scale parameter method: Fixed	value	▼ Value: 1.0 🗲
Covariance matrix: 💿 Mo	idel-based estimat	tor 🔘 Robust estimator
tterations	Output	
Singularity tolerance:		Descending Ollise data order
value order for categorical inputs:	 Ascenally (
OK 🕨 Run Cancel		Apply Reset

Figure 285. Choix des options expert

12. Exécutez le flux pour créer le nugget de modèle, qui est ajouté à l'espace de travail du flux et à la palette Modèles dans l'angle supérieur droit. Pour consulter les détails du modèle, cliquez avec le bouton droit de la souris sur le nugget et sélectionnez **Modifier** ou **Parcourir**.

Tests des effets du modèle

Tests of Model Effects

	Тур		
Source	Wald Chi-Square	df	Sig.
(Intercept)	.536	1	.464
Age in years	.358	1	.550
Duration of disease	.003	1	.958
Treatment group	.382	1	.537

Dependent Variable: Result

Model: (Intercept), Age in years, Duration of disease, Treatment group

Figure 286. Tests des effets pour le modèle Effets principaux

Aucun des effets du modèle n'est significatif d'un point de vue statistique ; toutefois, toute différence observable dans les effets du traitement a un intérêt du point de vue clinique. Nous allons donc ajuster un modèle réduit avec, pour seule caractéristique du modèle, le traitement.

Ajustement du modèle avec le traitement pour seule caractéristique

- 1. Dans l'onglet Champs du noeud Modèles linéaires généralisés, cliquez sur **Utiliser les paramètres personnalisés**.
- 2. Sélectionnez *result* comme cible.
- 3. Sélectionnez treatment comme seule entrée.

V Treatment-only	
	0 - 🗆
Fields Model Expert Analyze Annotation	8
O Use type node settings	Use custom settings
Target: 🔒 result	,
Inputs: \delta treatment	
	×
Partilizz	
	•
Splits:	
Use weight field	▼]
Target field represents number of events occ	urring in a set of trials
Trials field:	
Fixed value	
Number of trials: 10 🖨	
OK 🕨 Run Cancel	Apply Reset

Figure 287. Sélection des options de champ

4. Exécutez le flux et ouvrez le nugget de modèle résultant.

Dans le nugget de modèle, sélectionnez l'onglet Avancé et accédez au bas de la liste.

Estimations des paramètres

Parameter Estimates

			95% Wald Cont	fidence Interval	Hypoth	nesis Test	~
Parameter	В	Std. Error	Lower	Upper	Wald Chi-Square	df	Sig.
(Intercept)	-1.442	.5012	-2.425	460	8.282	1	.004
[Treatment group=1]	.378	.6288	855	1.610	.361	1	.548
[Treatment group=0]	0 ^a						
(Scale)	1 ^b						

Dependent Variable: Result

Model: (Intercept), Treatment group

a. Set to zero because this parameter is redundant.

b. Fixed at the displayed value.

Figure 288. Estimations des paramètres pour le modèle avec le traitement pour seule caractéristique

L'effet du traitement (la différence, pour le prédicteur linéaire, entre les deux niveaux de traitement ; c'est-à-dire le coefficient de [*treatment=1*]) n'est toujours pas significatif d'un point de vue statistique. Il suggère uniquement que le traitement *A* [*treatment=0*] semble meilleur que le traitement *B* [*treatment=1*] car l'estimation des paramètres pour le traitement *B* est supérieure à celle du traitement *A* et est donc associée à une probabilité de réapparition accrue dans les 12 premiers mois. Le prédicteur linéaire (constante + effet du traitement) est une estimation de log(-log(1-P(recur_{12,t})), où recur_{12,t}) est la probabilité de réapparition à 12 mois pour le traitement t(=*A* ou *B*). Ces probabilités prédites sont générées pour chaque observation du jeu de données.

Réapparition prédite et probabilités de survie

💟 Derive	×
	0.0
Derive as: Conditional	
Settings Annotations	
Mode: 💿 Single 🔘 Multiple	
Derive field:	
precur	
Derive as: Conditional T Field type: Y <default></default>	
If:	
Then:	
Else:	
OK Cancel	Apply Reset

Figure 289. Options des paramètres du noeud Dériver

- 1. Pour chaque patient, le modèle détermine le score du résultat prédit et de la probabilité de ce résultat prédit. De façon à visualiser les probabilités de réapparition prédites, copiez le modèle généré dans la palette et reliez un noeud Dériver.
- 2. Dans l'onglet Paramètres, saisissez le champ de calcul precur.
- 3. Choisissez de le calculer comme champ Conditionnel.
- 4. Cliquez sur le bouton représentant une calculatrice pour ouvrir le Générateur de formules pour la condition If (Si).

	and the second		atia		III. Ioido		
Return				Туре	Field -	Storage	
Boolean	-	<u> </u>	tem	8	result	Integer	
Boolean			mod	- 1	duration	Integer	
Boolean				8	treatment	Integer	
Boolean				A	age	Integer	
Boolean				8	\$G-result	Integer	
Boolean				Ì	\$GP-result	Real	
Boolean		and	or	1	\$GP-0	Real	
Boolean		notO		1	\$GP-1	Real	
Integer	-		<u> </u>	1	\$GRP-result	Real	
	Return Boolean Boolean Boolean Boolean Boolean Boolean Integer	Return Boolean Boolean Boolean Boolean Boolean Boolean Boolean Boolean Boolean Colean Boolean Boolean Colean Colea	Return Boolean Boolean Boolean Boolean Boolean Boolean Boolean Integer	Return * rem Boolean / mod Boolean > >= Boolean > >= Boolean > Boolean > Boolean > Boolean > Boolean > Boolean Boolean	Return * rem Type Boolean / mod //mod Boolean //mod //mod Boolean /mod //mod	Return * rem Type Field Boolean / mod area area area Boolean >>= area area area Boolean >>= area area area Boolean >>= area \$GP-result Boolean and or \$GP-1 \$GP-result Boolean area \$GP-result \$GP-result	Return Type Field Storage Boolean Integer Integer Boolean Integer Integer

Figure 290. Noeud Dériver : Générateur de formules pour la condition If

- 5. Insérez le champ *\$G-result* dans la formule.
- 6. Cliquez sur OK.

Le champ de calcul *precur* prendra la valeur de la formule **Then (Donc)** lorsque *\$G-result* est égal à 1 et la valeur de la formule **Else (Sinon)** lorsqu'il est égal à 0.

General Functions		-	+	**	1 3	Fields		
Function -	Return			uiv	Туре	Field	Storage	
_integer(ITEM)	Boolean	-	*	rem	8	result	Integer	-
_real(ITEM)	Boolean		4	mod		duration	Integer	
_number(ITEM)	Boolean				-	treatment	Integer	
_string(ITEM)	Boolean					age	Integer	
_date(ITEM)	Boolean				8	\$G-result	Integer	
time(ITEM)	Boolean		-		A	\$GP-result	Real	
timestamp(ITEM)	Boolean		and	or	A	\$GP-0	Real	
datetime(ITEM)	Boolean		not()		A	\$GP-1	Real	
_integer(ITEM)	Integer	+		5		\$GRP-result	Real	
								-

Figure 291. Noeud Dériver : Générateur de formules pour l'expression Then

- 7. Cliquez sur le bouton représentant une calculatrice pour ouvrir le Générateur de formules pour la formule **Then**.
- 8. Insérez le champ *\$GP-result* dans la formule.

9. Cliquez sur OK.

				. I. I. a. E	Fields		
Function -	Return			Туре	Field	Storage	
s_integer(ITEM)	Boolean	-	* rem	8	result	Integer	_
s_real(ITEM)	Boolean		mod		duration	Integer	
_number(ITEM)	Boolean			8	treatment	Integer	
_string(ITEM)	Boolean				age	Integer	
_date(ITEM)	Boolean			8	\$G-result	Integer	
_time(ITEM)	Boolean			A	\$GP-result	Real	
_timestamp(ITEM)	Boolean		and or	A	\$GP-0	Real	
_datetime(ITEM)	Boolean		not() ><	A	\$GP-1	Real	
_integer(ITEM)	Integer	-			\$GRP-result	Real	
: _integer(ITEM) sturns a value of true if I <u>C</u> heck expression bef	TEM type is an int ore saving	eger.	Otherwise, reti	urns a va	lue of false.		

Figure 292. Noeud Dériver : Générateur de formules pour l'expression Else

- **10**. Cliquez sur le bouton représentant une calculatrice pour ouvrir le Générateur de formules pour la formule **Else**.
- 11. Saisissez 1- dans la formule puis insérez-y le champ \$GP-result.
- 12. Cliquez sur OK.

😡 precur					×
				0	- 🗖
Derive as: Condition	al				
Settings Annotations					
Mod	e: 🔘 Sir	ngle 🔘 Multiple	э		
Derive field:					
precur					
Derive as: Conditional T Field type: If:	-				
'\$G-result'					
Then:					
'\$GP-result'					
Else:					
1-'\$GP-result'					
OK Cancel				Apply	Reset

Figure 293. Options des paramètres du noeud Calculer

13. Liez un noeud Table au noeud Dériver et exécutez-le.

🔟 Table	e (10 fi	ields, 4	3 records	s) #1	มีเรา	A 36		-0	X
Table	Annotati	ons	zenerale						
	result	duration	treatment	age	\$G-result	\$GP-result	\$GP-0	\$GP-1	
1	1	2	1	48	0	0.708	0.708	0.292	-
2	0	1	1	73	0	0.708	0.708	0.292	
3	0	1	1	54	0	0.708	0.708	0.292	
4	0	2	1	58	0	0.708	0.708	0.292	
5	0	1	0	56	0	0.789	0.789	0.211	
6	0	2	0	49	0	0.789	0.789	0.211	
7	0	1	1	71	0	0.708	0.708	0.292	
8	0	1	0	41	0	0.789	0.789	0.211	
9	0	1	1	23	0	0.708	0.708	0.292	
10	1	1	1	37	0	0.708	0.708	0.292	
11	0	1	1	38	0	0.708	0.708	0.292	
12	0	2	1	76	0	0.708	0.708	0.292	
13	0	2	0	38	0	0.789	0.789	0.211	
14	1	1	0	27	0	0.789	0.789	0.211	
15	1	1	1	47	0	0.708	0.708	0.292	
16	0	1	0	54	0	0.789	0.789	0.211	
17	1	1	1	38	0	0.708	0.708	0.292	
18	1	2	1	27	0	0.708	0.708	0.292	
19	0	2	0	58	0	0.789	0.789	0.211	
20	0	1	1	75	0	0.708	0.708	0.292	*
		and the second						E Contraction ()	Γ
								6	
								C.	JR

Figure 294. Probabilités prédites

La probabilité estimée que les patients voient réapparaître leur ulcère dans les 12 premiers mois est de 0,211 pour les patients recevant le traitement *A* et de 0,292 pour ceux qui reçoivent le traitement *B*. Notez que $1-P(\text{recur}_{12, t})$ est la probabilité de survie à 12 mois, laquelle peut s'avérer plus intéressante pour les analystes de la survie.

Modélisation de la probabilité de réapparition par période

Ce modèle présente un inconvénient : il ignore les informations recueillies lors du premier examen. En effet, chez de nombreux patients, l'ulcère n'est pas réapparu durant les six premiers mois. Pour obtenir un " meilleur " modèle, il serait nécessaire de modéliser une réponse binaire qui enregistre si l'événement est survenu ou n'est pas survenu durant chaque intervalle. L'ajustement de ce modèle nécessite une reconstruction du jeu de données d'origine, disponible dans le fichier *ulcer_recurrence_recoded.sav*. Ce fichier contient deux variables supplémentaires :

- Period, qui enregistre si l'observation correspond à la période du premier ou du second examen
- *Result by period,* qui enregistre si une réapparition est survenue ou non pour le patient donné durant la période donnée.

Chaque observation d'origine (le patient) fournit une observation pour chaque intervalle où il demeure dans l'ensemble de risques. Par exemple, le patient 1 fournit deux observations : une observation pour la période du premier examen durant laquelle aucune réapparition n'est survenue et une observation pour la période du second examen durant laquelle une réapparition a été enregistrée. Le patient 10, en revanche, ne fournit qu'une seule observation car une réapparition a été enregistrée dans la première période. Les patients 16, 28 et 34 ont abandonné l'étude après six mois et ne fournissent donc qu'une seule observation au nouveau jeu de données.

1. Ajoutez un noeud source Fichier Statistics pointant sur *ulcer_recurrence_recoded.sav* dans le dossier *Demos*.



Figure 295. Flux d'échantillons relatif à la prévision de la réapparition des ulcères

2. Dans l'onglet Filtrer du noeud source, excluez *id*, *time* et *result*.

Ulcer_recurrence_recode	ed.sav fresh currence_recoded.sa is	× •
7		Fields: 8 in, 3 filtered, 0 renamed, 5 out
Field -	Fitter	Field
id	×>	id
age	\rightarrow	age
duration	\rightarrow	duration
treatment	\rightarrow	treatment
time	— × →	time
result	→	result
period	\rightarrow	period
result2	\rightarrow	result2
View current fields View OK Cancel	unused field settings	Apply Reset

Figure 296. Filtrage des champs superflus

3. Dans l'onglet Types du noeud source, définissez le rôle du champ *result2* sur **Cible** et son niveau de mesure sur **Indicateur**. Le rôle de tous les autres champs doit être défini sur **Entrée**.

Field Measurement Values Missing Check Role age Continuous [23,76] None Input duration Ordinal 1,2 None Input treatment Nominal 0,1 None Input period Ordinal 1,2 None Input result2 Flag 1/0 None Input	Preview 2 Refresh \$CLEO_DEMOSAlcer_recurrence_recoded.sav Data Filter Types Annotations					
Field Measurement Values Missing Check Role age Continuous [23,76] None Input duration Ordinal 1,2 None Input treatment Nominal 0,1 None Input period Ordinal 1,2 None Input result2 Flag 1/0 None Input	√ - 00 0	Mead Value	es Clear	Values	Clear All Value	s
age Continuous [23,76] None Input duration Ordinal 1,2 None Input treatment Nominal 0,1 None Input period Ordinal 1,2 None Input result2 Flag 1/0 None Input	Field -	Measurement	Values	Missing	Check	Role
duration Ordinal 1,2 None Input treatment Nominal 0,1 None Input period Ordinal 1,2 None Input result2 Flag 1/0 None Target	> age	🔗 Continuous	[23,76]		None	🔪 Input
treatment Nominal 0,1 None Input period Ordinal 1,2 None Input result2 Flag 1/0 None Target	duration	- Ordinal	1,2		None	🔪 Input
period I Ordinal 1,2 None Input result2 I Flag 1/0 None I Target	> treatment	💑 Nominal	0,1		None	🔪 Input
>result2 🖁 Flag 1/0 None 🥥 Target	> period	📲 Ordinal	1,2		None	🔪 Input
	result2	🖁 Flag	1/0		None	🔘 Target
View current fields 💦 View upueed field eattings	View currect	fielde O View up 1999	I field actting			

Figure 297. Définition du rôle de champ

4. Ajoutez un noeud Re-trier et spécifiez *period*, *duration*, *treatment* et *age* comme ordre des entrées. Le fait d'avoir *period* comme première entrée (et de ne pas inclure la caractéristique de constante dans le modèle) permet d'ajuster un ensemble complet de variables factices pour capturer les effets de la période.

💟 Fiel	d Reorder		×
Reorde	Preview Apportations		0
O Custo	om Order	O Automatic Sort	
Type:	Name: 🔺 🔻 Storage:		
Туре	Field	Storage	
	[other fields]		
	period	🚫 Integer	
-	duration	🚫 Integer	
-	treatment	🚫 Integer	
	age	今 Integer	
			Ť
Clear L	Jnused		
Note: Fi	elds added down stream of this node are no	t reordered.	
ок	Cancel		Apply Reset

Figure 298. Réorganisation des champs afin qu'ils soient entrés dans le modèle de la façon souhaitée

5. Dans le noeud Modèles linéaires généralisés, cliquez sur l'onglet Modèle.

🙀 result2	
Fields Model Expert Analyze Annotations	
Model name: O Auto Custom	
☑ Use partitioned data	
Build model for each split	
Model type: Main effects only Main effects and all two-way interactions	
Offset:	
● Variable	
Offset field:	
Fixed value Value: 0.0	
Base category for flag target: First (Lowest)	
Include intercept in model	
OK Run Cancel	Apply Reset

Figure 299. Choix des options de modèle

- 6. Sélectionnez **Premiers (valeur la plus faible)** comme catégorie de référence pour la cible. Ce choix indique que la seconde catégorie est l'événement d'intérêt et que son effet sur le modèle se situe dans l'interprétation des estimations des paramètres.
- 7. Désélectionnez l'option Inclure une constante dans le modèle.
- 8. Cliquez sur l'onglet Expert et sélectionnez Expert pour activer les options de modélisation expert.

💟 result	
	0
Fields Model Expert Analyze Annotations	
Mode: 🔘 Simple 💿 Expert	
Target Field Distribution and Link Function	
The distribution that you choose determines which lin	ik functions are available.
Distribution: Binomial	Parameters
	Parameter for negative binomial:
	Specify value Value: 1.0 🗧
	© Estimate
	Parameter for Tweedle:
Link function: Complementary log-log	Power: 0.0
Method and iteration settings are not available if Distrib	ution = Normal and Link
Function = Identity.	
Method:	Maximum Fisher scoring iterations:
Scale parameter method: Fixed value	Value:
Covariance matrix: O Model-based estimator	r 🔘 Robust estimator
tterations Output	
Singularity tolerance: 1E-007 T	
Value order for categorical inputs: O Ascending ()	/Descending 🔘 Use data order
OK 🕨 Run Cancel	Apply Reset

Figure 300. Choix des options expert

- 9. Sélectionnez **Binomial** pour la distribution et **Log-log complémentaire** pour la fonction de lien.
- 10. Sélectionnez **Valeur fixe** comme méthode d'estimation du paramètre d'échelle et conservez la valeur par défaut 1,0.
- 11. Sélectionnez l'ordre des catégories des facteurs **Décroissant**. Cela signifie que la première catégorie de chaque facteur constitue la catégorie de référence. L'effet de cette sélection sur le modèle porte sur l'interprétation des estimations de paramètre.
- **12.** Exécutez le flux pour créer le nugget de modèle, qui est ajouté à l'espace de travail du flux et à la palette Modèles dans l'angle supérieur droit. Pour consulter les détails du modèle, cliquez avec le bouton droit de la souris sur le nugget et sélectionnez **Modifier** ou **Parcourir**.

Tests des effets du modèle

Tests of Model Effects

	Туре III					
Source	Wald Chi-Square	df	Sig.			
Period	.464	1	.496			
Age in years	.314	1	.575			
Duration of disease	.000	1	.988			
Treatment group	.117	1	.732			

Dependent Variable: Result by period

Model: Period, Age in years, Duration of disease, Treatment group

Figure 301. Tests des effets pour le modèle Effets principaux

Aucun des effets du modèle n'est significatif d'un point de vue statistique ; toutefois, toute différence observable dans les effets de la période et du traitement a un intérêt du point de vue clinique ; nous allons donc ajuster un modèle réduit, avec ces seules caractéristiques de modèle.

Ajustement du modèle réduit

- 1. Dans l'onglet Champs du noeud Modèles linéaires généralisés, cliquez sur **Utiliser les paramètres personnalisés**.
- 2. Sélectionnez *result2* comme cible.
- 3. Sélectionnez *period* et *treatment* comme entrées.

6			
😡 Period-Treatment			X
			0
Fields Model Expert A	nalyze Annotations		
O Use type node settings		Ose custom settings	
Target: 🔒 result2			J
Inputs: 📑 period			×
Partition:			
Splits:			×
Use weight field			-1
 Target field represents nu Variable 	mber of events occurri	ng in a set of trials	
Trials field:			-
Sixed value			
Number of trials:	10 🜲		
OK 🕨 Run Car	ncel		Apply Reset

Figure 302. Sélection des options de champ

4. Exécutez le noeud et parcourez le modèle généré. Ensuite, copiez le modèle généré dans la palette, reliez un noeud Table et exécutez-le.

Estimations des paramètres

Parameter Estimates

			95% Wald Conf	fidence Interval	Hypoth	nesis Test	
Parameter	В	Std. Error	Lower	Upper	Wald Chi-Square	df	Sig.
[Period=2]	-1.794	.5792	-2.929	659	9.597	1	.002
[Period=1]	-2.206	.5912	-3.365	-1.047	13.926	1	.000
[Treatment group=1]	.195	.6279	-1.035	1.426	.097	1	.756
[Treatment group=0]	0 ^a	2.44					
(Scale)	1 ^b						

Dependent Variable: Result by period

Model: Period, Treatment group

a. Set to zero because this parameter is redundant.

b. Fixed at the displayed value.

Figure 303. Estimations des paramètres pour le modèle avec le traitement pour seule caractéristique

L'effet du traitement n'est toujours pas significatif d'un point de vue statistique. Il suggère uniquement que le traitement *A* semble meilleur que le traitement *B* car l'estimation des paramètres pour le traitement *B* est associée à une probabilité de réapparition accrue durant les 12 premiers mois. Les valeurs de période sont, d'un point de vue statistique, significativement différentes de 0, mais ceci est dû au fait qu'un terme de constante n'est pas ajusté. L'effet de la période (la différence entre les prédicteurs linéaires de [*period=1*] et [*period=2*]) n'est pas significatif, d'un point de vue statistique, comme l'indiquent les tests des effets du modèle. Le prédicteur linéaire (effet de la période + effet du traitement) est une estimation de log($-\log(1-P(\text{recur}_{p, t}))$, où P($\text{recur}_{p, t}$) est la probabilité de réapparition au cours de la période p(=1 ou 2, représentant six mois ou 12 mois), compte tenu du traitement t(=*A* ou *B*). Ces probabilités prédites sont générées pour chaque observation du jeu de données.

Réapparition prédite et probabilités de survie

💟 Derive	×
	0.0
Derive as: Conditional	
Settings Annotations	
Mode: 💿 Single 🔘 Multiple	
Derive field:	
precur	
Derive as: Conditional T Field type: Y <default></default>	
If:	
Then:	
Else:	
OK Cancel	Apply Reset

Figure 304. Options des paramètres du noeud Dériver

- 1. Pour chaque patient, le modèle détermine le score du résultat prédit et de la probabilité de ce résultat prédit. De façon à visualiser les probabilités de réapparition prédites, copiez le modèle généré dans la palette et reliez un noeud Dériver.
- 2. Dans l'onglet Paramètres, saisissez le champ de calcul precur.
- 3. Choisissez de le calculer comme champ Conditionnel.
- 4. Cliquez sur le bouton représentant une calculatrice pour ouvrir le Générateur de formules pour la condition If (Si).

General Functions		-			Fields		
Function -	Return			Тур	e Field -	Storage	
_integer(ITEM)	Boolean	4	L.		result2	Integer	
_real(ITEM)	Boolean			mod	period	Integer	
_number(ITEM)	Boolean			>=	duration	Integer	_
_string(ITEM)	Boolean		\mathbf{a}	7	treatment	Integer	_
_date(ITEM)	Boolean				age	Integer	
_time(ITEM)	Boolean				\$G-result2	Integer	
_timestamp(ITEM)	Boolean		and	or 🧳	\$GP-result2	Real	
_datetime(ITEM)	Boolean		notO		\$GP-0	Real	
_integer(ITEM)	Integer	-			\$GP-1	Real	
_addedime((TEM)	Integer	-			\$GP-1	Real	

Figure 305. Noeud Dériver : Générateur de formules pour la condition If

- 5. Insérez le champ *\$G-result2* dans la formule.
- 6. Cliquez sur OK.

Le champ de calcul *precur* prendra la valeur de la formule **Then** lorsque *\$G-result2* est égal à 1 et la valeur de la formule **Else** lorsqu'il est égal à 0.

		1	dika		Fields	
Function -	Return			Туре	Field -	Storage
_integer(ITEM)	Boolean	4	* rem	8	result2	Integer
_real(ITEM)	Boolean				period	Integer
_number(ITEM)	Boolean		>>=		duration	Integer
_string(ITEM)	Boolean		6 6	1 🗸 🗌	treatment	Integer
_date(ITEM)	Boolean				age	Integer
_time(ITEM)	Boolean			1 8	\$G-result2	Integer
_timestamp(ITEM)	Boolean		and or		\$GP-result2	Real
_datetime(ITEM)	Boolean		not() ><	1	\$GP-0	Real
integer(ITEM)	Integer				\$GP-1	Real

Figure 306. Noeud Dériver : Générateur de formules pour l'expression Then

- 7. Cliquez sur le bouton représentant une calculatrice pour ouvrir le Générateur de formules pour la formule **Then**.
- 8. Insérez le champ *\$GP-result2* dans la formule.

9. Cliquez sur OK.

General Functions		-	+ **	个皆	Fields		-
Function -	Return			Туре	Field	Storage	
s_integer(ITEM)	Boolean	-	× rem	8	result2	Integer	-
_real(ITEM)	Boolean		mod		period	Integer	
_number(ITEM)	Boolean		>>=		duration	Integer	
_string(ITEM)	Boolean		00	i 👼	treatment	Integer	
_date(ITEM)	Boolean				age	Integer	
_time(ITEM)	Boolean				\$G-result2	Integer	
_timestamp(ITEM)	Boolean		and or		\$GP-result2	Real	
_datetime(ITEM)	Boolean		not() ><		\$GP-0	Real	
_integer(ITEM)	Integer	-			\$GP-1	Real	-
	1						

Figure 307. Noeud Dériver : Générateur de formules pour l'expression Else

- **10**. Cliquez sur le bouton représentant une calculatrice pour ouvrir le Générateur de formules pour la formule **Else**.
- 11. Saisissez 1- dans la formule puis insérez-y le champ \$GP-result2.
- 12. Cliquez sur OK.

🔽 precur	X
Derive as: Conditional	
Settings Annotations	
Mode: 💿 Single 🔘 Multiple	
Derive field:	
precur	
Derive as: Conditional T Field type: 🖋 <default> T</default>	
'\$G-result2'	
Then:	
'\$GP-result2'	
Else:	
1-'\$GP-result2'	
OK Cancel	Apply Reset

Figure 308. Options des paramètres du noeud Dériver

13. Liez un noeud Table au noeud Dériver et exécutez-le.

III Table (11 fields, 78 records) #3									
じ <u>F</u> ile	📄 Edi	t 🏷	<u>G</u> enerate		B	14			(🖉 🔪
Table Annotations									
	result2	period	duration	treatment	age	\$G-result2	\$GP-result2	\$GP-0	\$GP-1
1	0	1	2	1	48	0	0.875	0.875	0.125
2	1	2	2	1	48	0	0.817	0.817	0.183
3	0	1	1	1	73	0	0.875	0.875	0.125
4	0	2	1	1	73	0	0.817	0.817	0.183
5	0	1	1	1	54	0	0.875	0.875	0.125
6	0	2	1	1	54	0	0.817	0.817	0.183
7	0	1	2	1	58	0	0.875	0.875	0.125
8	0	2	2	1	58	0	0.817	0.817	0.183
9	0	1	1	0	56	0	0.896	0.896	0.104
10	0	2	1	0	56	0	0.847	0.847	0.153
11	0	1	2	0	49	0	0.896	0.896	0.104
12	0	2	2	0	49	0	0.847	0.847	0.153
13	0	1	1	1	71	0	0.875	0.875	0.125
14	0	2	1	1	71	0	0.817	0.817	0.183
15	0	1	1	0	41	0	0.896	0.896	0.104
16	0	2	1	0	41	0	0.847	0.847	0.153
17	0	1	1	1	23	0	0.875	0.875	0.125
18	0	2	1	1	23	0	0.817	0.817	0.183
19	1	1	1	1	37	0	0.875	0.875	0.125
20	0	1	1	1	38	0	0.875	0.875	0.125
	4	5011		NEGHTSON					1
									6

Figure 309. Probabilités prédites

Tableau 3. Probabilités de réapparition estimées

Traitement	6 mois	12 mois
А	0,104	0,153
В	0,125	0,183

A partir de ces probabilités de réapparition estimées, la probabilité de survie sur 12 mois peut être estimée sous la forme $1-(P(\text{recur}_{1, t}) + P(\text{recur}_{2, t}) \times (1-P(\text{recur}_{1, t})))$. Par conséquent, pour chaque traitement :

A: 1 - (0,104 + 0,153*0,896) = 0,759

 $B: 1 - (0,125 + 0,183^*0,875) = 0,715$

ce qui montre de nouveau une préférence, significative d'un point de vue autre que statistique, pour *A* comme étant le meilleur traitement.

Récapitulatif

A l'aide des modèles linéaires généralisés, vous avez ajusté une série de modèles de régression log-log complémentaires pour des données de survie avec censure par intervalle. Même si le choix du traitement *A* est privilégié, l'obtention d'un résultat significatif du point de vue statistique peut nécessiter une étude plus importante. Toutefois, d'autres pistes sont à explorer avec les données existantes.

• Il peut être utile de réajuster le modèle avec des effets d'interaction, notamment entre *Period* et *Treatment group*.

Vous trouverez des explications sur le fondement mathématique des méthodes de modélisation utilisées dans IBM SPSS Modeler dans le *Guide des algorithmes IBM SPSS Modeler*.

Procédures apparentées

La procédure Modèles linéaires généralisés constitue un outil puissant qui permet d'adapter toutes sortes de modèles.

- La procédure Equations d'estimation généralisées étend le modèle linéaire généralisé pour autoriser les mesures répétées.
- La procédure Modèles mixtes linéaires permet d'adapter des modèles pour des variables d'échelle dépendantes avec un composant aléatoire et/ou des mesure répétées.

Lectures recommandées

Reportez-vous aux documents suivants pour plus d'informations sur les modèles linéaires généralisés :

Cameron, A. C. et P. K. Trivedi. 1998. *Regression Analysis of Count Data*. Cambridge : Cambridge University Press. Dobson, A. J. 2002. *An Introduction to Generalized Linear Models*, 2 ed. Boca Raton, FL: Chapman & Hall/CRC. Hardin, J. W. et J. M. Hilbe. 2003. *Generalized Linear Models and Extension*. Station, TX: Stata Press. McCullagh, P. et J. A. Nelder. 1989. *Generalized Linear Models*, 2nd ed. Londres : Chapman & Hall.

Chapitre 23. Utilisation de la régression de Poisson pour analyser les taux de dommage aux navires (modèles linéaires généralisés)

Un modèle linéaire généralisé peut être utilisé pour ajuster une régression de Poisson pour l'analyse des données d'effectif. Par exemple, un jeu de données présenté et analysé ailleurs (²) relate les dommages que les vagues causent aux cargos. Le nombre d'incidents peut être modélisé comme se produisant selon la fréquence définie par un test Poisson en fonction des prédicteurs, et le modèle résultant peut permettre de déterminer les types de navire les plus exposés aux dommages.

Cet exemple utilise le flux *ships_genlin.str*, qui fait référence au fichier de données *ships.sav*. Le fichier de données se trouve dans le dossier *Demos* et le fichier de flux dans le sous-dossier *streams*.

La modélisation des calculs des cellules brutes peut se révéler insatisfaisante dans cette situation car le *total des mois de service* varie selon le type de navire. Les variables qui mesurent le degré d''' exposition '' au risque sont traitées dans le modèle linéaire généralisé en tant que variables de décalage. D'autre part, une régression de Poisson considère que le log de la variable dépendante est linéaire dans les prédicteurs. Par conséquent, pour utiliser les modèles linéaires généralisés pour ajuster une régression de Poisson aux taux d'accidents, vous devez utiliser le *logarithme du total des mois de service*.

Ajustement d'une régression de Poisson "surdispersée"

1. Ajoutez un noeud source Statistics pointant vers *ships.sav* dans le dossier *Demos*.



Figure 310. Flux d'échantillons utilisé pour l'analyse des taux de dommage

^{2.} McCullagh, P. et J. A. Nelder. 1989. Generalized Linear Models, 2nd ed. Londres : Chapman & Hall.

[©] Copyright IBM Corp. 1994, 2018

2. Dans l'onglet Filtrer du noeud source, excluez le champ *months_service*. Les valeurs transformées en log de cette variable sont contenues dans *log_months_service*, qui sera utilisé dans l'analyse.

🕜 ships.sav		
Preview Refr	resh	0
\$CLEO_DEMOS/ships.sav	r -	
Data Filter Types Annotations		
		Fields: 6 in, 1 filtered, 0 renamed, 5 out
Field -	Filter	Field
type	\rightarrow	type
construction	\rightarrow	construction
operation	\rightarrow	operation
months_service	★ →	months_service
log_months_service	\rightarrow	log_months_service
damage_incidents	\rightarrow	damage_incidents
View current fields O View (unused field settings	
OK Cancel		Apply Reset

Figure 311. Filtrage d'un champ inutile

(Vous pouvez également régler le rôle de ce champ sur **Aucun** dans l'onglet Types au lieu de l'exclure ou sélectionner les champs que vous souhaitez utiliser dans le noeud de modélisation.)

- **3**. Dans l'onglet Types du noeud source, définissez le rôle du champ *damage_incidents* sur **Cible**. Le rôle de tous les autres champs doit être défini sur **Entrée**.
- 4. Cliquez sur Lire les valeurs pour instancier les données.

🕜 ships.sav					
	eview) 😰 Refresh)			0
\$CLEO_	DEMOS/ships.sav				
Data Filter	Annotations				
4. 00 0	💌 🚺 🕨 Read Val	ues Clear V	Values	Clear All Valu	ies
Field	Measurement	Values	Missing	Check	Role
🚫 type 👘	💑 Nominal	1,2,3,4,5		None	🔪 Input
🔆 construction 💧	📶 Ordinal	60,65,70,75		None	🔪 Input
🔷 operation 👘	📶 Ordinal	60,75		None	🔪 Input
🛞 log_months	🔗 Continuous	[3.806662		None	🛇 None
🚫 damage_inc	🔗 Continuous	[0,58]		None	🔘 Target
View current f	ields 🔘 View unus	ed field settings			
OK Cancel	ieius 🥥 view unus	ea neia settings			Apply Reset

Figure 312. Définition du rôle de champ

- 5. Reliez un noeud Modèles linéaires généralisés au noeud source. Dans le noeud Modèles linéaires généralisés, cliquez sur l'onglet **Modèle**.
- 6. Sélectionnez *log_months_service* comme variable de décalage.

😡 Over dispersed Poisson	\mathbf{X}
	0
Fields Model Expert Analyze Annotations	
Model name: O Auto O Custom Overdispersed Poisson	
☑ Use partitioned data	
☑ Build model for each split	
Model type: O Main effects only O Main effects and all two-way interactions	
Offset:	
Offset field: 🔗 log_months_service	-
© Fixed value Value: 0.0 €	
Base category for flag target: Last (Highest) 🔽	
Include intercept in model	
OK Run Cancel	Apply Reset

Figure 313. Choix des options de modèle

7. Cliquez sur l'onglet **Expert** et sélectionnez **Expert** pour activer les options de modélisation expert.

😡 Overdispersed Poisson	
Fields Model Expert Analyze Annotations	
Mode: O Simple O Expert	
Target Field Distribution and Link Function	
The distribution that you choose determines which li	nk functions are available.
Distribution: Poisson	Parameters
	Parameter for negative binomial:
	Specify value Value: 1.0 🖨
	Estimate
	Parameter for Tweedie: 1.5
Link function: Log	Power: 0.0
Method and iteration settings are not available if Distrik Function = Identity. ┌─Parameter Estimation	oution = Normal and Link
Method:	Maximum Fisher scoring iterations:
Scale parameter method: Pearson Chi-square	Value:
Covariance matrix: O Model-based estimato	r 🔘 Robust estimator
Iterations	
Singularity tolerance:	
Value order for categorical inputs: O Ascending 🧕) Descending 🔘 Use data order
OK 🕨 Run Cancel	Apply Reset

Figure 314. Choix des options expert

- 8. Sélectionnez Poisson comme distribution pour la réponse et Log comme fonction de lien.
- 9. Sélectionnez **Khi-carré de Pearson** comme méthode d'estimation du paramètre d'échelle. Le paramètre d'échelle est généralement considéré comme étant égal à 1 dans une régression de Poisson. Cependant, McCullagh et Nelder utilisent l'estimation khi-deux de Pearson pour obtenir des estimations de variance et des niveaux de signification plus prudents.

- 10. Sélectionnez l'ordre des catégories des facteurs **Décroissant**. Cela signifie que la première catégorie de chaque facteur constitue la catégorie de référence. L'effet de cette sélection sur le modèle porte sur l'interprétation des estimations de paramètre.
- 11. Cliquez sur **Exécuter** pour créer le nugget de modèle qui est ajouté à l'espace de travail du flux et à la palette Modèles en haut à droite. Pour consulter les détails du modèle, cliquez avec le bouton droit de la souris sur le nugget et sélectionnez **Modifier** ou **Parcourir**, puis sur l'onglet **Avancé**.

	Value	df	Value/df
Deviance	38.695	25	1.548
Scaled Deviance	22.883	25	
Pearson Chi-Square	42.275	25	1.691
Scaled Pearson Chi-Square	25.000	25	
Log Likelihoodª	-68.281		
Akaike's Information Criterion (AIC)	154:562		
Finite Sample Corrected AIC (AICC)	162.062		
Bayesian Information Criterion (BIC)	168.299		
Consistent AIC (CAIC)	177.299		

Statistiques de qualité d'ajustement

Dependent Variable: Number of damage incidents

Model: (Intercept), type, construction, operation, offset = log_months_ service

- a. The full log likelihood function is displayed and used in computing information criteria.
- b. Information criteria are in small-is-better form.

Figure 315. Statistiques de qualité d'ajustement

Le tableau des statistiques de qualité d'ajustement contient des mesures permettant de comparer les modèles en concurrence. Par ailleurs, la *valeur/df* de la déviance et des statistiques du khi-carré de Pearson donne des estimations pour le paramètre d'échelle. Ces valeurs doivent être proches de 1,0 pour une régression de Poisson. Des valeurs supérieures à 1,0 indiquent que l'ajustement du modèle surdispersé peut être judicieux.

Test composite

Omnibus Test^a

Likelihood Ratio Chi-Square	df	Sig.
63.650	8	.000

Dependent Variable: Number of damage incidents

Model: (Intercept), Year of construction, Period of operation, Ship type, offset = Logarithm of aggregate months of service

a. Compares the fitted model against the intercept-only model.

Figure 316. Test composite

Le test composite est un test du khi-carré du rapport de vraisemblance du modèle actuel par rapport au modèle nul (dans ce cas, la constante). Une valeur de signification inférieure à 0,05 indique que le modèle actuel permet d'obtenir de meilleurs résultats que le modèle nul.

Tests des effets du modèle

Tests of Model Effects

	Туре III				
Source	Wald Chi-Square	df	Sig.		
(Intercept)	2138.657	1	.000		
Year of construction	17.242	3	.001		
Period of operation	6.249	1	.012		
Ship type	15.415	4	.004		

Dependent Variable: Number of damage incidents

Model: (Intercept), Year of construction, Period of operation, Ship type, offset = Logarithm of aggregate months of service

Figure 317. Tests des effets du modèle

Chaque terme du modèle est testé afin de déterminer s'il présente un effet. Les termes dont les valeurs de signification sont inférieures à 0,05 ont un effet visible. Chacun des termes des effets principaux contribue au modèle.

Estimations des paramètres

Parameter Estimates

			95% Wald Confidence Interval		Hypothesis Test		
Parameter	В	Std. Error	Lower	Upper	Wald Chi-Square	df	Sig.
(Intercept)	-6.406	.2828	-6.960	-5.852	513.238	1	.000
[Year of construction=75]	.453	.3032	141	1.048	2.236	1	.135
[Year of construction=70]	.818	.2208	.386	1.251	13.743	1	.000
[Year of construction=65]	.697	.1946	.316	1.079	12.835	1	.000
[Year of construction=60]	0 ^a						2.00
[Period of operation=75]	.384	.1538	.083	.686	6.249	1	.012
[Period of operation=60]	0 ^a			2			~
[Ship type=5]	.326	.3067	276	.927	1.127	1	.288
[Ship type=4]	076	.3779	817	.665	.040	1	.841
[Ship type=3]	687	.4279	-1.526	.151	2.581	1	.108
[Ship type=2]	543	.2309	996	091	5.536	1	.019
[Ship type=1]	0 ^a					÷	~
(Scale)	1.691 ^b						

Dependent Variable: Number of damage incidents

Model: (Intercept), Year of construction, Period of operation, Ship type, offset = Logarithm of aggregate months of service

a. Set to zero because this parameter is redundant.

b. Computed based on the Pearson chi-square.

Figure 318. Estimations des paramètres

Le tableau des estimations de paramètre récapitule l'effet de chaque prédicteur. Bien que l'interprétation des coefficients dans ce modèle soit difficile de par la nature de la fonction de lien, les signes des coefficients des covariables et les valeurs relatives des coefficients des niveaux de facteur peuvent fournir des informations importantes sur les effets des prédicteurs dans le modèle.

- Pour les covariables, les coefficients positifs (négatifs) indiquent les relations positives (inverses) entre les prédicteurs et la sortie. Une valeur croissante d'une covariable avec un coefficient positif correspond à un taux croissant d'incidents provoquant des dommages.
- Dans le cas des facteurs, un niveau de facteur présentant un coefficient supérieur indique un impact plus important des dommages. Le signe d'un coefficient d'un niveau de facteur dépend de l'effet de ce niveau de facteur par rapport à la modalité de référence.

Vous pouvez tirer les conclusions suivantes en fonction des estimations des paramètres :

- Le type de navire *B* [*type=2*] présente un taux de dommage considérablement inférieur du point de vue statistique (valeur *p* de 0,019) (coefficient estimé de -0,543) que le type *A* [*type=1*], la catégorie de référence. Le type *C* [*type=3*] présente en réalité un paramètre estimé inférieur à celui de *B*, mais la variabilité de l'estimation de *C* trouble l'effet. Reportez-vous aux moyennes marginales estimées pour toutes les relations entre les niveaux du facteur.
- Les navires construits entre 1965 et 1969 [*construction=65*] et 1970 et 1974 [*construction=70*] présentent des taux de dommage considérablement supérieurs du point de vue statistique (valeurs *p* <0,001)

(coefficients estimés de 0,697 et de 0,818, respectivement) que ceux construits entre 1960 et 1964 *[construction=60]*, la catégorie de référence. Reportez-vous aux moyennes marginales estimées pour toutes les relations entre les niveaux du facteur.

• Les navires en service entre 1975 et 1979 [*operation*=75] présentent des taux de dommage considérablement supérieurs du point de vue statistique (valeur *p* de 0,012) (coefficient estimé de 0,384) que ceux en service entre 1960 et 1974 [*operation*=60].

Ajustement des modèles alternatifs

Le problème concernant la régression de Poisson "surdispersée" est qu'il n'existe pas de manière formelle de la tester par rapport à la régression de Poisson "standard". Toutefois, un test formel conseillé afin de déterminer l'existence d'une surdispersion consiste à effectuer un test de rapport de vraisemblance entre une régression de Poisson "standard" et une régression binomiale négative, l'ensemble des autres paramètres étant égaux. En cas d'absence de surdispersion dans la régression de Poisson, la statistique 2x(log de vraisemblance du modèle de Poisson - log de vraisemblance du modèle binomial négatif) doit présenter une proportion de mélange : la moitié de sa masse de probabilité sur 0 et le reste dans une distribution Khi-deux avec 1 degré de liberté.

1. Sélectionnez Valeur fixe comme méthode d'estimation du paramètre d'échelle. Par défaut, cette valeur est 1.

Negative Binomia	ι	$\overline{\mathbf{X}}$
Fields Model Expert	Analyze Annotations	
Mode: O Simple O Expe	ert	
Target Field Distribution a	nd Link Function	
The distribution that you c	hoose determines which lini	k functions are available.
Distribution: Negative	binomial 🔨	Parameters
		Parameter for negative binomial:
		ـ ◎ Specify value Value: 1.0 ≑
		© Estimate
		Parameter for Tweedie:
Link function: Log	T	Power: 0.0
Method and iteration setting	gs are not available if Distribu	ution = Normal and Link
Function = Identity.		
Parameter Estimation		
Method:	Hybrid	Maximum Fisher scoring iterations:
Scale parameter method:	Fixed value	Value: 1.0 🖨
Covariance matrix:	Model-based estimator	◎ Robust estimator
tterations	Output	
Singularity tolerance:		
Value order for categorica	l inputs: 🔘 Ascending 🍥	Descending 🔘 Use data order
OK 🕨 Run	Cancel	Apply

Figure 319. Onglet Expert

- 2. Pour ajuster la régression binomiale négative, copiez et collez le noeud Modèles linéaires généralisés, liez-le au noeud source, ouvrez le nouveau noeud, puis cliquez sur l'onglet **Expert**.
- **3**. Sélectionnez la **loi binomiale négative**. Conservez la valeur par défaut de 1 pour le paramètre secondaire.
- 4. Exécutez le flux et accédez à l'onglet Avancé dans les nuggets de modèle nouvellement créés.

Statistiques de qualité d'ajustement

	Value	df	Value/df
Deviance	38.695	25	1.548
Scaled Deviance	38.695	25	
Pearson Chi-Square	42.275	25	1.691
Scaled Pearson Chi-Square	42.275	25	
Log Likelihood ^a	-68.281		
Akaike's Information Criterion (AIC)	154.562		
Finite Sample Corrected AIC (AICC)	162.062		
Bayesian Information Criterion (BIC)	168.299		
Consistent AIC (CAIC)	177.299		

Dependent Variable: Number of damage incidents

Model: (intercept), type, construction, operation, offset = log_months_ service

 The full log likelihood function is displayed and used in computing information criteria.

b. Information criteria are in small-is-better form.

Figure 320. Statistiques de qualité d'ajustement pour la régression de Poisson standard

Le log de vraisemblance indiqué pour la régression de Poisson standard est -68.281. Comparez ce chiffre à celui du modèle binomial négatif.

	Value	df	Value/df
Deviance	11.145	25	.446
Scaled Deviance	11.145	25	
Pearson Chi-Square	8.815	25	.353
Scaled Pearson Chi-Square	8.815	25	
Log Likelihoodª	-83.725		
Akaike's Information Criterion (AIC)	185.450		
Finite Sample Corrected AIC (AICC)	192.950		
Bayesian Information Criterion (BIC)	199.187		
Consistent AIC (CAIC)	208.187		

Dependent Variable: Number of damage incidents

Model: (Intercept), type, construction, operation, offset = log_months_ service

a. The full log likelihood function is displayed and used in

computing information criteria.

b. Information criteria are in small-is-better form.

Figure 321. Statistiques de qualité d'ajustement pour la régression binomiale négative

Le log de vraisemblance indiqué pour la régression binomiale négative est –83.725. Il est en réalité *inférieur* au log de vraisemblance de la régression de Poisson, ce qui indique (sans le test de rapport de vraisemblance) que cette régression binomiale négative n'offre pas d'amélioration par rapport à la régression de Poisson.
Toutefois, la valeur de 1 choisie pour le paramètre secondaire de la loi binomiale négative peut ne pas être optimale pour ce jeu de données. Une autre manière de tester la surdispersion consiste à ajuster un modèle binomial négatif avec un paramètre secondaire égal à 0 et à demander le test du multiplicateur de Lagrange dans la boîte de dialogue Sortie de l'onglet Expert. Si le test n'est pas concluant, la surdispersion ne doit pas poser de problème à ce jeu de données.

Récapitulatif

A l'aide des modèles linéaires généralisés, vous avez ajusté trois modèles différents pour les données d'effectif. Il s'est avéré que la régression binomiale négative n'offre aucune amélioration par rapport à la régression de Poisson. La régression de Poisson surdispersée semble offrir une alternative raisonnable au modèle de Poisson standard, mais il n'existe pas de test formel permettant de choisir l'un plutôt que l'autre.

Vous trouverez des explications sur le fondement mathématique des méthodes de modélisation utilisées dans IBM SPSS Modeler dans le *Guide des algorithmes IBM SPSS Modeler*.

Procédures apparentées

La procédure Modèles linéaires généralisés constitue un outil puissant qui permet d'adapter toutes sortes de modèles.

- La procédure Equations d'estimation généralisées étend le modèle linéaire généralisé pour autoriser les mesures répétées.
- La procédure Modèles mixtes linéaires permet d'adapter des modèles pour des variables d'échelle dépendantes avec un composant aléatoire et/ou des mesure répétées.

Lectures recommandées

Reportez-vous aux documents suivants pour plus d'informations sur les modèles linéaires généralisés :

Cameron, A. C. et P. K. Trivedi. 1998. *Regression Analysis of Count Data*. Cambridge : Cambridge University Press. Dobson, A. J. 2002. *An Introduction to Generalized Linear Models*, 2 ed. Boca Raton, FL: Chapman & Hall/CRC. Hardin, J. W. et J. M. Hilbe. 2003. *Generalized Linear Models and Extension*. Station, TX: Stata Press. McCullagh, P. et J. A. Nelder. 1989. *Generalized Linear Models*, 2nd ed. Londres : Chapman & Hall.

Chapitre 24. Ajustement d'une régression gamma à des déclarations de sinistre automobile (modèles linéaires généralisés)

Un modèle linéaire généralisé permet d'ajuster une régression gamma pour l'analyse de données d'intervalle positif. Par exemple, un jeu de données présenté et analysé ailleurs (³) porte sur les déclarations de sinistre automobile. La somme moyenne des déclarations peut être modélisée comme suivant une distribution gamma, en utilisant une fonction de lien inverse pour associer la moyenne de la variable dépendante à une combinaison linéaire de prédicteurs. Afin de représenter le nombre variable de déclarations servant à calculer les sommes moyennes des déclarations, vous pouvez indiquer la pondération de mise à l'échelle *Number of claims*.

Cet exemple utilise le flux *car-insurance_genlin.str*, qui fait référence au fichier de données *car_insurance_claims.sav*. Le fichier de données se trouve dans le dossier *Demos* et le fichier de flux dans le sous-dossier *streams*.

Création du flux

1. Ajoutez un noeud source Fichier Statistics pointant sur *car_insurance_claims.sav* dans le dossier *Demos.*



Figure 322. Flux d'échantillons relatif à la prévision des déclarations de sinistre automobile

- 2. Dans l'onglet Types du noeud source, définissez le rôle du champ *claimamt* sur **Cible**. Le rôle de tous les autres champs doit être défini sur **Entrée**.
- 3. Cliquez sur Lire les valeurs pour instancier les données.

^{3.} McCullagh, P. et J. A. Nelder. 1989. Generalized Linear Models, 2nd ed. Londres : Chapman & Hall.

Data Fitter Typ	es Annotations	lues Clear '	Values	Clear All Valu	es
Field -	Measurement	Values	Missing	Check	Role
> holderage 🛛	👖 Ordinal	1,2,3,4,5,		None	🔪 Input
> vehiclegroup 💧	Ы Nominal	1,2,3,4		None	🔪 Input
> vehicleage 🔡	📶 Ordinal	1,2,3,4		None	🔪 Input
🕻 claimamt 🛛	🔗 Continuous	[11,850]		None	🔘 Target
🔪 nclaims 🛛 💡	🔗 Continuous	[0,434]		None	O None
View current fi	elds 🔘 View unus	ed field settings	3		

Figure 323. Définition du rôle de champ

- 4. Reliez un noeud Modèles linéaires généralisés au noeud source. Dans le noeud Modèles linéaires généralisés, cliquez sur l'onglet Champs.
- 5. Sélectionnez le champ de pondération *nclaims*.

🔽 claimamt			X
			0-0
Fields Model Expert	Analyze Annotations		
OUse type node setting:	s	O Use custom settings	
Target:			_
Inputs:			×
Partition:			-
Splits:			×
🗹 Use weight field	🔗 nclaims		-1
Target field represents	number of events occurri	ng in a set of trials	
Trials field:			-
Fixed value Number of trials:	10 🗘		
OK 🕨 Run	Cancel		Apply Reset

Figure 324. Sélection des options de champ

6. Cliquez sur l'onglet Expert et sélectionnez Expert pour activer les options de modélisation expert.

💟 claimamt	×
Fields Model Expert Analyze Annotat	ions
Mode: 🔘 Simple 🔘 Expert	
Target Field Distribution and Link Function-	
The distribution that you choose determines w	/hich link functions are available.
Distribution: Gamma	Parameters
	Parameter for negative binomial:
	Specify value Value: 1.0 🖨
	Estimate
	Parameter for Tweedle:
Link function: Power	Power: 1.0
Method and iteration settings are not available it	f Distribution = Normal and Link
Function = Identity.	
Parameter Estimation	
Method: Hybrid	Maximum Fisher scoring iterations: 1
Scale parameter method: Pearson Chi-squar	re 🔻 Value: 1.0 荣
Covariance matrix: Model-based es	stimator 🔘 Robust estimator
Iterations	1
Singularity tolerance: 1E-007 T	
Value order for categorical inputs: O Ascend	ling 💿 Descending 🔘 Use data order
OK 🕨 Run Cancel	Apply

Figure 325. Choix des options expert

- 7. Sélectionnez la distribution de réponse Gamma.
- 8. Sélectionnez la fonction de lien **Puissance**, puis saisissez l'exposant de la fonction de puissance -1,0. Il s'agit d'un lien inverse.
- **9**. Sélectionnez la méthode d'estimation du paramètre d'échelle **Khi-carré de Pearson**. Il s'agit de la méthode utilisée par McCullagh et Nelder. Nous allons donc la suivre de manière à reproduire leurs résultats.
- 10. Sélectionnez l'ordre des catégories des facteurs **Décroissant**. Cela signifie que la première catégorie de chaque facteur constitue la catégorie de référence. L'effet de cette sélection sur le modèle porte sur l'interprétation des estimations de paramètre.
- 11. Cliquez sur **Exécuter** pour créer le nugget de modèle qui est ajouté à l'espace de travail du flux et à la palette Modèles en haut à droite. Pour afficher les détails du modèle, cliquez avec le bouton droit de la souris sur le nugget de modèle et choisissez **Modifier** ou **Parcourir**, puis sélectionnez l'onglet Avancé.

Estimations des paramètres

Parameter Estimates

			95% Wald Confidence Interval		Hypothesis Test		
Parameter	В	Std. Error	Lower	Upper	Wald Chi-Square	df	Sig.
(Intercept)	.003	.0004	.003	.004	66.593	1	.000
[Policyholder age=8]	.001	.0004	.000	.002	4.898	1	.027
[Policyholder age=7]	.001	.0004	.000	.002	5.046	1	.025
[Policyholder age=6]	.001	.0004	.000	.002	5.740	1	.017
[Policyholder age=5]	.001	.0004	.001	.002	10.682	1	.001
[Policyholder age=4]	.000	.0004	.000	.001	1.268	1	.260
[Policyholder age=3]	.000	.0004	.000	.001	.720	1	.396
[Policyholder age=2]	.000	.0004	001	.001	.054	1	.816
[Policyholder age=1]	0 ^a	8				8	8
[Vehicle age=4]	.004	.0004	.003	.005	88.175	1	.000
[Vehicle age=3]	.002	.0002	.001	.002	53.013	1	.000
[Vehicle age=2]	.000	.0001	.000	.001	13.191	1	.000
[Vehicle age=1]	0 ^a	8				3	8
[Vehicle group=4]	001	.0002	002	001	61.883	1	.000
[Vehicle group=3]	001	.0002	001	.000	13.039	1	.000
[Vehicle group=2]	3.765E-5	.0002	.000	.000	.050	1	.823
[Vehicle group=1]	0 ^a	8					8
(Scale)	1.209 ^b						

Dependent Variable: Average cost of claims

Model: (Intercept), Policyholder age, Vehicle age, Vehicle group

- a. Set to zero because this parameter is redundant.
- b. Computed based on the Pearson chi-square.

Figure 326. Estimations des paramètres

Le test composite et les tests d'effets de modèle (non affichés) indiquent que le modèle permet d'obtenir de meilleurs résultats que le modèle nul et que chaque caractéristique effet principal contribue à ce modèle. Le tableau des estimations de paramètre contient les mêmes valeurs que celles obtenues par McCullagh et Nelder pour les niveaux de facteur et le paramètre d'échelle.

Récapitulatif

Grâce aux modèles linéaires généralisés, vous venez d'ajuster une régression gamma aux données concernant les déclarations. Même si une fonction de lien canonique pour la distribution gamma a été utilisée dans ce modèle, un lien log produit également des résultats raisonnables. En général, il est difficile, voire impossible, de comparer directement des modèles avec différentes fonctions de lien. Toutefois, le lien log constitue un cas particulier du lien de puissance, où l'exposant est égal à 0. Par

conséquent, vous pouvez comparer les déviances d'un modèle avec un lien log et un modèle avec un lien de puissance pour déterminer celui qui offre le meilleur ajustement (reportez-vous, par exemple, à la section 11.3 concernant McCullagh et Nelder).

Vous trouverez des explications sur le fondement mathématique des méthodes de modélisation utilisées dans IBM SPSS Modeler dans le *Guide des algorithmes IBM SPSS Modeler*.

Procédures apparentées

La procédure Modèles linéaires généralisés constitue un outil puissant qui permet d'adapter toutes sortes de modèles.

- La procédure Equations d'estimation généralisées étend le modèle linéaire généralisé pour autoriser les mesures répétées.
- La procédure Modèles mixtes linéaires permet d'adapter des modèles pour des variables d'échelle dépendantes avec un composant aléatoire et/ou des mesure répétées.

Lectures recommandées

Reportez-vous aux documents suivants pour plus d'informations sur les modèles linéaires généralisés :

Cameron, A. C. et P. K. Trivedi. 1998. *Regression Analysis of Count Data*. Cambridge : Cambridge University Press. Dobson, A. J. 2002. *An Introduction to Generalized Linear Models*, 2 ed. Boca Raton, FL: Chapman & Hall/CRC. Hardin, J. W. et J. M. Hilbe. 2003. *Generalized Linear Models and Extension*. Station, TX: Stata Press. McCullagh, P. et J. A. Nelder. 1989. *Generalized Linear Models*, 2nd ed. Londres : Chapman & Hall.

Chapitre 25. Classification des échantillons de cellules (SVM)

Support Vector Machine (SVM) est une technique de classification et de régression particulièrement adaptée aux larges jeux de données. Un large jeu de données est un ensemble contenant un nombre important de prédicteurs, comme c'est le cas dans le domaine de la bio-informatique (l'application des technologies de l'information aux données biochimiques et biologiques).

Un chercheur en médecine a obtenu un jeu de données contenant les caractéristiques d'un certain nombre d'échantillons de cellules humaines supposées favoriser le développement du cancer. L'analyse des données originales indiquait que de nombreuses caractéristiques différaient considérablement entre les échantillons bénins et malins. Ce chercheur en médecine souhaite développer un modèle SVM qui peut utiliser les valeurs des caractéristiques de ces cellules dans des échantillons d'autres patients pour savoir au plus tôt si leurs échantillons peuvent être bénins ou malins.

Cet exemple utilise le flux nommé *svm_cancer.str*, disponible dans le dossier *Demos* du sous-dossier des *flux*. Le fichier de données est *cell_samples.data*. Pour plus d'informations, voir la rubrique «Dossier Demos», à la page 4.

Cet exemple utilise un jeu de données disponible au public dans le référentiel d'apprentissage automatique. Ce jeu de données est constitué de plusieurs centaines d'enregistrements d'échantillons de cellules humaines, chacun d'entre eux contenant les valeurs d'un ensemble de caractéristiques des cellules. Les champs de chaque enregistrement sont :

Nom du champ	Description
ID	Identifiant du patient
Clump	Epaisseur de l'agglutination
UnifSize	Uniformité de la taille des cellules
UnifShape	Uniformité de la forme des cellules
MargAdh	Adhésion marginale
SingEpiSize	Taille des cellules épithéliales
BareNuc	Noyau nu
BlandChrom	Chromatine terne
NormNucl	Nucléole normal
Mit	Mitoses
Class	Bénigne ou maligne

Dans cet exemple, nous utilisons un jeu de données contenant un nombre relativement petit de prédicteurs dans chaque enregistrement.

Création du flux



Figure 327. Flux d'échantillons présentant la modélisation SVM

1. Créez un nouveau flux et ajoutez un noeud source Délimité pointant vers *cell_samples.data* dans le dossier *Demos* de votre installation IBM SPSS Modeler.

Examinons les données du fichier source.

- 2. Ajoutez un noeud Table au flux.
- 3. Liez le noeud Table au noeud Délimité et exécutez le flux.

🛙 Table (11 fields, 699 records) 📃 🗆 🔽										
じ <u>F</u> ile	📄 Ed	it 👋 Ger	herate						0	×
Table 4	Annotatio	ons								
	hifSize	UnifShape	MargAdh	SingEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class	
1		1	1	2	1	3	1	1	2	4
2		4	5	7	10	3	2	1	2	Ľ
3		1	1	2	2	3	1	1	2	
4		8	1	3	4	3	7	1	2	
5		1	3	2	1	3	1	1	2	
6	1	10	8	7	10	9	7	1	4	
7		1	1	2	10	3	1	1	2	
8		2	1	2	1	3	1	1	2	
9		1	1	2	1	1	1	5	2	
10		1	1	2	1	2	1	1	2	
11		1	1	1	1	3	1	1	2	
12		1	1	2	1	2	1	1	2	
13		3	3	2	3	4	4	1	4	
14		1	1	2	3	3	1	1	2	
15		5	10	7	9	5	5	4	4	
16		6	4	6	1	4	3	1	4	
17		1	1	2	1	2	1	1	2	
18		1	1	2	1	3	1	1	2	
19		7	6	4	10	4	1	2	4	
20		1	1	2	1	3	1	1	2	
	4	-	و و واد دی از ایش	A CONTRACTOR OF A CONTRACT			Long-to-to-to-to-to-to-to-to-to-to-to-to-to-	in the second		Г

Figure 328. Données source de SVM

Le champ *ID* contient les identifiants du patient. Les caractéristiques des échantillons de cellules de chaque patient se trouvent dans les champs *Clump* à *Mit*. Les valeurs vont de 1 à 10, 1 étant le plus proche de bénin.

Le champ *Class* contient le diagnostique, confirmé par plusieurs procédures médicales, établissant si les échantillons sont bénins (valeur = 2) ou malins (valeur = 4).

Type	Annotations				× •
4. 00	Read Val	ues Clear	Values	Clear All Va	alues
Field -	Measurement	Values	Missing	Check	Role
VUNITSIZE	🞸 conunuous	[1,10]		NONE	🔳 iniput 🖌
📿 UnifShape 👘	🞸 Continuous	[1,10]		None	🔪 Input 📂
🚫 MargAdh	🔗 Continuous	[1,10]		None	🔪 Input
🚫 SingEpiSize	🔗 Continuous	[1,10]		None	🔪 Input
A BareNuc	💑 Nominal	"1","10","		None	🔪 Input
SlandChrom	🖉 Continuous	[1,10]		None	> Input
NormNucl	Continuous	[1,10]		None	> Input
🔆 Mit	Continuous	[1.10]		None	hinput
🔆 Class	🎖 Flag	4/2		None	🔘 Target 🗖
OK Cance	fields 🔘 View unus	ed field settin	gs		Apply Reset

Figure 329. Paramètres du noeud type

4. Ajoutez un noeud type et liez-le au noeud Délimité.

5. Ouvrez le noeud type.

Nous voudrions que le modèle prédise la valeur de *Class* (c'est-à-dire, bénigne (=2) ou maligne (=4)). Ce champ ne pouvant avoir qu'une des deux valeurs possibles, il est nécessaire de modifier son niveau de mesure pour refléter ceci.

- 6. Dans la colonne **Mesure** du champ *Class* (le dernier de la liste), cliquez sur la valeur **Continu** et modifiez-la en **Indicateur**.
- 7. Cliquez sur Lire les valeurs.
- 8. Dans la colonne **Rôle**, définissez le rôle du champ *ID* (l'identifiant du patient) sur **Aucun**, cette valeur n'étant pas utilisée comme prédicteur ou comme cible du modèle.
- 9. Définissez le rôle de la cible, *Class*, sur **Cible** et laissez le rôle de tous les autres champs (prédicteurs) sur **Entrée**.
- 10. Cliquez sur OK.

Le noeud SVM propose plusieurs fonctions du noyau permettant son exécution. Comme il n'est pas évident de savoir quelle fonction est la plus appropriée à un jeu de données spécifique, nous choisirons plusieurs fonctions afin de comparer leurs résultats. Commençons par la fonction par défaut, RBF (Fonction radiale de base).

😡 c lass	-rbf					
						0
Fields	Model	Expert	Analyze	Annotations		
Model nar	ne:		() At	to 🔘 Custom	class-rbf	
🔽 Use p	artition	ed data				
🛃 Build r	nodel f	or each s	plit			
To select	fields	manually,	choose "U	se custom settir	ngs" on the Fields t	ab
Partition	n I					-1
Splits:						×
ОК		Run	Cancel			Apply Reset

Figure 330. Paramètres de l'onglet Modèle

- 11. Dans la palette Modélisation, reliez un noeud SVM au noeud type.
- 12. Ouvrez le noeud SVM. Dans l'onglet **Modèle**, cliquez sur l'option **Personnalisé** pour le **nom du modèle** et entrez *class-rbf* dans le champ de texte adjacent.

😡 class-rbf			
			0
Fields Model Expert Analy	yze Annotations		
Mode:	🔘 Simple 🧿 Exper	t	
📃 Append all probabilities (valid	only for categorical	targets)	
Stopping criteria:	1.0E-3 🔻		
Regularization parameter (C):	10 ≑		
Regression precision (epsilon):	0.1 ≑		
Kernel type:	RBF 💌		
RBF gamma:	0.1 🚔	Bias:	0 ≑
Gamma:	1 🗢	Degree:	3 🖨
OK 🕨 Run Cancel			Apply Reset

Figure 331. Onglet Expert - Paramètres par défaut

13. Dans l'onglet **Expert**, définissez le **Mode** sur **Expert** pour faciliter la lecture mais ne modifiez aucune des options par défaut. Veuillez noter que **type noyau** est défini sur **RBF** par défaut. Toutes les options sont grisées en mode Simple.

🖓 Clas	s					
Fields	Model	Expert	Analyze	Annotations		
[Model	Evaluatio	on				
Ca	alculate v	/ariable in	portance			
Proper	nsity Sco	ores (valid	l only for fl	ag targets)—		
Ca	alculate r	aw prope	ensity scor	es		
🔳 Ca	alculate a	adjusted p	ropensity	scores		
Based	lon			0	esting partition	n ⊚ ∀alidation partition
OK		Run	Cancel			Apply Reset

Figure 332. Paramètres de l'onglet Analyser

- 14. Dans l'onglet Analyser, sélectionnez la case Calculer l'importance de la variable.
- **15**. Cliquez sur **Exécuter**. Le nugget de modèle est placé dans le flux et dans la palette Modèles en haut à droite de l'écran.
- 16. Double-cliquez sur le nugget de modèle dans le flux.

Analyse des données



Figure 333. Graphique de l'importance des prédicteurs

Dans l'onglet Modèle, le graphique d'importance des prédicteurs présente l'effet relatif des différents champs sur la prévision. Ceci nous indique que *BareNuc* a bien l'effet le plus important alors que *UnifShape* et *Clump* sont également relativement importants.

- 1. Cliquez sur OK.
- 2. Liez un noeud Table au nugget de modèle class-rbf.
- 3. Ouvrez le noeud Table, puis cliquez sur Exécuter.

👜 Table	Edit	(us , 099 () <u>G</u> er	nerate					
Table	Annotation	s						
	gEpiSize	BareNuc	BlandChrom	NormNucl	Mit	Class	\$S-Class	\$SP-Class
1		1	3	1	1	2	2	0.992
2		10	3	2	1	2	4	0.899
3		2	3	1	1	2	2	0.994
4		4	3	7	1	2	4	0.915
5		1	3	1	1	2	2	0.992
6		10	9	7	1	4	4	0.999
7		10	3	1	1	2	2	0.907
8		1	3	1	1	2	2	0.997
9		1	1	1	5	2	2	0.997
10		1	2	1	1	2	2	0.996
11		1	3	1	1	2	2	0.999
12		1	2	1	1	2	2	0.999
13		3	4	4	1	4	2	0.514
14		3	3	1	1	2	2	0.989
15		9	5	5	4	4	4	0.991
16		1	4	3	1	4	4	0.691
17		1	2	1	1	2	2	0.997
18		1	3	1	1	2	2	0.995
19		10	4	1	2	4	4	0.996
20		1	3	1	1	2	2	0.986
	4						الحصار النجار	

Figure 334. Champs ajoutés pour la valeur de prévision et de confiance

4. Le modèle a créé deux champs supplémentaires. Faites défiler les sorties de la table vers la droite pour les voir :

Nouveau nom de champ	Description
\$S-Class	Valeur de la <i>classe</i> prédite par le modèle.
\$SP-Class	Score de propension de cette prévision (la probabilité qu'a cette prévision d'être vraie, une valeur de 0,0 à 1,0).

D'un coup d'oeil à la table, nous pouvons voir que les scores de propension (dans la colonne *\$SP-Class*) de la majorité des enregistrements sont assez élevés.

Mais il existe des exceptions importantes ; par exemple, l'enregistrement pour le patient 1041801 à la ligne 13 dont la valeur de 0,514 est anormalement basse. De plus, si l'on compare *Class* à *\$S-Class*, il est évident que ce modèle a effectué plusieurs prévisions incorrectes, même lorsque le score de propension était relativement élevé (par exemple, lignes 2 et 4).

Voyons si le résultat peut être meilleur en choisissant un autre type de fonction.

Essai d'une autre fonction

😡 class-pol	y				X
					0
Fields Mode	Expert	Analyze	Annotations		
Model name:		O Au	.to 💿 Custom	class-poly	
👿 Use partitio	ned data				
V Build mode	for each s	plit			
To select field	s manually,	choose "U	se custom setti	ngs" on the Fields ta	ıb
Partition:					-1
Splits:					×
ок	Run	ancel			Apply Reset

Figure 335. Définition d'un nouveau nom pour le modèle

- 1. Fermez la fenêtre de sortie Table.
- 2. Reliez un deuxième noeud de modélisation SVM au noeud type.
- 3. Ouvrez le nouveau noeud SVM.
- 4. Dans l'onglet **Modèle**, choisissez Personnalisé et saisissez *class-poly* comme nom de modèle.

😡 class-poly			
			0
Fields Model Expert Ana	alyze Annotations		
Mode:	🔘 Simple 🧿 Exper	t	
📃 Append all probabilities (vali	d only for categorical	targets)	
Stopping criteria:	1.0E-3 💌		
Regularization parameter (C):	10 ≑		
Regression precision (epsilon):	0.1 荣		
Kernel type:	Polynomial 🔻		
RBF gamma:	0.1 ≑	Bias:	0 🖨
Gamma:	1 ≑	Degree:	3 🗲
OK 🕨 Run Canc			Apply Reset

Figure 336. Paramètres de l'onglet Expert pour polynomial

5. Dans l'onglet Expert, définissez le Mode sur Expert.

- **6**. Définissez le **type du noyau** sur **Polynomial** et cliquez sur **Exécuter**. Le nugget de modèle *class-poly* est ajouté au flux et dans la palette Modèles en haut à droite de l'écran.
- 7. Connectez le nugget de modèle *class-rbf* au nugget de modèle *class-poly* (choisissez **Remplacer** dans la boîte de dialogue d'avertissement).
- 8. Liez un noeud Table au nugget *class-poly*.
- 9. Ouvrez le noeud Table, puis cliquez sur Exécuter.

Comparaison des résultats

違 <u>F</u> ile	📄 Edit	ų) Gene	erate [0
Table Annotations							
	ormNucl	Mit	Class	\$S-Class	\$SP-Class	\$S1-Class	\$SP1-Class
78	1	1	2	2	0.992	2	0.998
79		1	2	2	0.968	2	0.967
80		1	2	2	0.998	2	0.994
81		1	2	2	0.986	2	0.991
82		1	2	2	0.996	2	0.997
83		1	2	2	0.991	2	0.998
84	1	1	2	2	0.970	2	0.998
85	þ	7	4	4	0.992	4	1.000
86)	10	4	4	0.974	4	1.000
87		1	4	4	0.786	4	0.958
88	1	3	4	4	0.988	4	0.935
89	1	1	2	2	0.995	2	0.997
90	1	1	2	2	0.998	2	0.991
91		1	2	2	0.999	2	0.993
92		1	2	2	0.998	2	0.996
93		1	2	2	0.995	2	0.997
94		1	2	2	0.999	2	0.994
95		1	2	2	0.998	2	0.995
96		1	2	2	0.999	2	0.993
97		1	2	2	0.999	2	0.995
	4				(annotation)		

Figure 337. Champs ajoutés pour la fonction Polynomiale

1. Faites défiler les sorties de la table vers la droite pour voir les champs ajoutés.

Les champs générés pour le type de fonction Polynomiale sont appelés \$S1-Class et \$SP1-Class.

Les résultats pour Polynomiale semblent bien meilleurs. La majorité des scores de propension sont de 0,995 ou plus ce qui est très encourageant.

2. Pour confirmer l'amélioration dans le modèle, liez un noeud Analyse au nugget de modèle *class-poly*.

Ouvrez le noeud Analyse, puis cliquez sur Exécuter.

🔍 Analysis of	f [Class]				_ 🗆 🔀
😺 Eile 🏼 🖻 E	Edit 🐻	B 14			0 ×
Analysis Anr	notations		¢.		
Collapse A	II 🤤 Exp	and All			
Results for Individua	output field Cla al Models nparing \$S-Cla	ass ass with C	lass		
	Correct	684	97.85%	5	
	Wrong	15	2.15%	5	
	Total	699			
- Con	nparing \$S1-C	lass with	Class	_	
	Correct	699	100%		
	Wrong	0	0%		
	Total	699			
Agreem	ent between	\$S-Class	\$S1-Class		
Ag	ree	684	97.85%		
Dis	sagree	15	2.15%		
То	tal	699			
-Con	nparing Agree	ment with	Class		
	Correct	684	100%		
	Wrong	0	0%		
	Total	684			
					<u>ok</u>

Figure 338. Noeud Analyse

Cette technique qui emploie le noeud Analyse vous permet de comparer au moins deux nuggets de modèle de même type. La sortie du noeud Analyse indique que la fonction RBF prédit correctement 97,85% des cas, ce qui est relativement bon. Mais cette sortie indique que la fonction Polynomiale a correctement prédit le diagnostique pour chacun des cas. En pratique, il est peu probable que vous soyez confronté à une exactitude de 100 %. Vous pouvez néanmoins utiliser le noeud Analyse pour déterminer si le modèle a une exactitude acceptable pour votre application.

En fait, aucun des autres types de fonction (Sigmoïdale ou Linéaire) n'a donné de résultats aussi bons que ceux de la fonction Polynomiale sur ce jeu de données précis. Cependant, avec un autre jeu de données, les résultats pourraient facilement être différents. Par conséquent, il est utile d'essayer toutes les options disponibles.

Récapitulatif

Vous avez utilisé différents types de fonctions du noyau SVM afin de prédire une classification à partir de plusieurs attributs. Vous avez vu la façon dont différents noyaux donnent différents résultats pour le même jeu de données et comment mesurer l'amélioration d'un modèle par rapport à un autre.

Chapitre 26. Utilisation de la régression de Cox pour modéliser la durée jusqu'à l'attrition de la clientèle

Dans ses efforts pour réduire l'attrition de la clientèle, une entreprise de télécommunications s'intéresse à la modélisation de la "durée d'attrition" afin de déterminer les facteurs associés aux clients qui changent rapidement de service. A cette fin, un échantillon aléatoire de clients est sélectionné et la période pendant laquelle ils ont été client, qu'ils soient encore des clients actifs ou non, ainsi que différents autres champs sont extraits de la base de données.

Cet exemple utilise le flux *telco_coxreg.str*, qui fait référence au fichier de données *telco.sav*. Le fichier de données se trouve dans le dossier *Demos* et le fichier de flux dans le sous-dossier *streams*. Pour plus d'informations, voir la rubrique «Dossier Demos», à la page 4.

Création d'un modèle adapté

1. Ajoutez un noeud source de fichier Statistics pointant vers telco.sav dans le dossier Demos.



Figure 339. Flux d'échantillons pour l'analyse de la durée jusqu'à l'attrition

2. Dans l'onglet Filtrer du noeud source, excluez les champs *region, income, longten* de *wireten*, et *loglong* à*logwire*.

CLEO_DEMOSAtelco.s	Refresh	0	
Data Fitter Types Annotati	ons		
Field -	Filter	Fields: 42 in, 12 filtered, 0 rename Field	d, 30 ou
region	- × >	region	-
tenure	\rightarrow	tenure	
age	\rightarrow	age	
marital	\rightarrow	marital	
address	\rightarrow	address	
ncome	× →	income	
ed	\rightarrow	ed	
employ	\rightarrow	employ	
retire	\rightarrow	retire	
gender	\rightarrow	gender	-
View current fields Vie OK Cancel	w unused field setting	Apply	Reset

Figure 340. Filtrage des champs inutiles

(Vous pouvez également régler le rôle de ces champs sur **Aucun** dans l'onglet Types au lieu de l'exclure ou sélectionner les champs que vous souhaitez utiliser dans le noeud de modélisation.)

- **3**. Dans l'onglet Types du noeud source, définissez le rôle du champ *attrition* sur **Cible** et son niveau de mesure sur **Indicateur**. Le rôle de tous les autres champs doit être défini sur **Entrée**.
- 4. Cliquez sur Lire les valeurs pour instancier les données.

	Preview) 😰 Refresh D_DEMOS <i>i</i> telco.sav	1			
ata Filter 1	Types Annotations				
4 - 00	🗪 🜗 Read Va	alues Clea	ar Values	Clear All V	alues
Field -	Measurement	Values	Missina	Check	Role
pager	💑 Nominal	0,1		None	🔪 Input
internet	💑 Nominal	0,1		None	🔪 Input
> callid	💑 Nominal	0,1		None	🔪 Input
🔉 callwait	💑 Nominal	0,1		None	🔪 Input
👌 forward	💑 Nominal	0,1		None	🔪 Input
🔉 confer	💑 Nominal	0,1		None	🔪 Input
ebill	💑 Nominal	0,1		None	🔪 Input
lninc 👔	🔗 Continuous	[2.19722		None	🔪 Input
🕻 custcat	💑 Nominal	1,2,3,4		None	🔪 Input
🕻 churn	🎖 Flag	1/0		None	🔘 Target
View currer	nt fields 🔘 View unu:	sed field settir	igs		

Figure 341. Définition du rôle de champ

5. Liez un noeud Cox au noeud source : dans l'onglet **Champs**, sélectionnez la variable de durée de survie *tenure*.

😡 churi	n					
COX						0
Fields	Model	Expert	Settings	Annotations		
Survival tin	ne: 🥖	> tenure				_]
🔘 Use ty	pe noc	le setting:	5	0	Use custom settings	
Target:						-
Inputs:						×
Partition:						-1
Splits:						×
ОК		Run	Cancel			oply <u>R</u> eset

Figure 342. Sélection des options de champ

- 6. Cliquez sur l'onglet Modèle.
- 7. Sélectionnez la méthode de sélection de variables Pas à pas.

😡 churn					
COX					0
Fields Mode	Expert	Settings	Annotations		
Model name:		() A	uto 🔘 Custom		
👿 Use partitio	oned data				
👿 Build mode	l for each s	plit			
Method:	Stepwise				
Groups:					
Model type:	🔘 Main ef	ffects 🔘 (Custom		
Model terms:					
					××
ОК	Run	Cancel		(Apply Reset

Figure 343. Choix des options de modèle

8. Cliquez sur l'onglet Expert et sélectionnez Expert pour activer les options de modélisation expert.

9. Cliquez sur Résultat.

isplay:	At each step 🔘 At last step
Cl for exp(B)	Correlation of estimates
Display baseline function	
ots	
Survival 📝 Hazard	E Log minus log 🔲 One minus survival
ot a separate line for each va	ilue:
alue to use for plots:	
/alue to use for plots: Field	Value ~
/alue to use for plots: Field If tenure	Value
/alue to use for plots: Field ∲ tenure ∳ age	Value
/alue to use for plots: Field	Value
✓alue to use for plots: Field ✓ tenure ✓ age ✓ marital ✓ address	Value ← Mean Mean Mean Mean
Value to use for plots: Field ✓ tenure ✓ age marital ✓ address ✓ ed	Value ← Mean Mean Mean Mean Mean
/alue to use for plots: Field ∳ tenure ∳ age marital address ed ed ∳ employ	Value ← Mean Mean Mean Mean Mean Mean
/alue to use for plots: Field ✓ tenure ✓ age marital ✓ address ed ✓ employ ✓ retire	Value
Value to use for plots: Field ✓ tenure ✓ age marital ✓ address ed ed ✓ employ ✓ retire ✓ gender	Value ← Mean Mean

Figure 344. Choix des options de sortie avancées

- 10. Sélectionnez Survie et Risque comme tracés à créer puis cliquez sur OK.
- 11. Cliquez sur Exécuter pour créer le nugget de modèle qui est ajouté au flux et à la palette Modèles en haut à droite. Pour en afficher les détails, double-cliquez sur le nugget dans le flux. Pour commencer, examinez l'onglet Sorties Avancées.

Observations censurées

		N	Percent
Cases available in analysis	Event ^a	274	27.4%
	Censored	726	72.6%
	Total	1000	100.0%
Cases dropped	Cases with missing values	0	0.0%
	Cases with negative time	0	0.0%
	Censored cases before the earliest event in a stratum	0	0.0%
	Total	0	0.0%
Total		1000	100.0%

Figure 345. Récapitulatif du traitement des observations

La variable de statut identifie si l'événement s'est produit pour une observation donnée. Si l'événement ne s'est pas produit, l'observation est dite censurée. Les observations censurées ne sont pas utilisées pour le calcul des coefficients de régression mais sont utilisées pour calculer le risque de référence. Le récapitulatif du traitement des observations affiche 726 observations censurées. Il s'agit des clients qui ne sont pas partis.

Codages de variables catégorielles

	-	Frequency	(1) ^b	(2)	(3)	(4)
marital ^a	0=Unmarried	505	1			
	1=Married	495	0			
edª	1=Did not complete high school	204	1	0	0	0
	2=High school degree	287	0	1	0	0
	3=Some college	209	0	0	1	0
	4=College degree	234	0	0	0	1
6	5=Post-undergraduate degree	66	0	0	0	0
retire ^a	.00=No	953	1			
gondorð	1.00=Yes	47	0	1		
gender ^a	O=Male	483	1			
	1=Female	517	0			
tollfree ^a	0=No	526	1)		
6	1=Yes	474	0			
equip ^a	0=No	614	1			
	1=Yes	386	0)		
callcard ^a	0=No	322	1			
	1=Yes	678	0			
wireless ^a	0=No	704	1	1		
	1=Yes	296	0			
multline ^a	0=No	525	1			
	1=Yes	475	0			
voiceª	0=No	696	1			
	1=Yes	304	0			
multline ^a voice ^a pager ^a	0=No	739	1			
	1=Yes	261	0			
internet ^a	0=No	632	1			
3	1=Yes	368	0		[]	
callid ^a	0=No	519	1			
	1=Yes	481	0		i i	
callwait ^a	0=No	515	1			
	1=Yes	485	0			
forward ^a	0=No	507	1		i i	
0	1=Yes	493	0			
confer ^a	0=No	498	1			
	1=Yes	502	0	1		
ebill ^a	0=No	629	1			
	1=Yes	371	0			
custcat ^a	1=Basic service	266	1	0	0	
	2=E-service	217	0	1	0	
	3=Plus service	281	0	0	1	
	4=Total service	236	0	0	0	

Figure 346. Codages de variables catégorielles

Les codages de variables catégorielles sont une référence utile pour l'interprétation des coefficients de régression des covariables catégorielles, et particulièrement des variables dichotomiques. Par défaut, la catégorie de référence est la "dernière" catégorie. Ainsi, par exemple, bien que les clients dont le statut est *Marié* aient des valeurs de variable de 1 dans le fichier de données, ils sont codés comme 0 pour les besoins de la régression.

Sélection des variables

		0	/erall (score)	Change (From Previo	us Step	Change F	rom Previo	us Block
Step	-2 Log Likelihoo d	Chi- square	df	Siq.	Chi- square	df	Siq.	Chi- square	df	Siq.
1 ^a	3392.536	162.303	1	.000	133.828	1	.000	133.828	1	.000
2 ^b	3087.314	249.392	2	.000	305.222	1	.000	439.050	2	.000
3°	3027.085	328.426	3	.000	60.229	1	.000	499.279	3	.000
4 ^d	2990.790	347.197	4	.000	36.294	1	.000	535.574	4	.000
5 ^e	2973.790	362.673	5	.000	17.000	1	.000	552.574	5	.000
6 ^f	2958.796	376.140	6	.000	14.994	1	.000	567.568	6	.000
79	2945.503	384.717	7	.000	13.293	1	.000	580.861	7	.000
8h	2936.993	417.341	8	.000	8.510	1	.004	589.371	8	.000
9 ⁱ	2926.000	423.911	9	.000	10.994	1	.001	600.364	9	.000
10 ^j	2917.551	428.078	10	.000	8.449	1	.004	608.813	10	.000
11 ^k	2913.308	436.837	11	.000	4.243	1	.039	613.056	11	.000
12 ¹	2908.078	440.158	12	.000	5.230	1	.022	618.286	12	.000

2908.078 440.158 12 .000 5.230 1 .022 6
 a. Variable(s) Entered at Step Number 1: callcard
 b. Variable(s) Entered at Step Number 2: longmon
 c. Variable(s) Entered at Step Number 3: equip
 d. Variable(s) Entered at Step Number 4: employ
 e. Variable(s) Entered at Step Number 6: multime
 f. Variable(s) Entered at Step Number 6: woice
 g. Variable(s) Entered at Step Number 7: address
 h. Variable(s) Entered at Step Number 7: address
 h. Variable(s) Entered at Step Number 9: ebill
 j. Variable(s) Entered at Step Number 9: ebill
 j. Variable(s) Entered at Step Number 10: callid
 k. Variable(s) Entered at Step Number 11: internet
 l. Variable(s) Entered at Step Number 12: reside
 m. Beginning Block Number 0, initial Log Likelihood function: -2 Log likelihood: 3526.364
 n. Beginning Block Number 1. Method = Forward Stepwise (Likelihood Ratio)

Figure 347. Tests composites

Le processus de création de modèle utilise un algorithme Pas à pas ascendant. Les tests composites sont des mesures de la performance des modèles. La modification du khi-carré de l'étape précédente est la différence entre le log de vraisemblance -2 du modèle à l'étape précédente et à l'étape actuelle. Si cette étape devait ajouter une variable, l'inclusion serait logique si la signification de cette modification était inférieure à 0,05. Si cette étape devait supprimer une variable, l'exclusion serait logique si la signification de cette modification était supérieure à 0,10. Au cours des douze étapes, douze variables sont ajoutées au modèle.

		В	SE	Wald	df	Siq.	Exp(B)
Step 12	address	035	.009	14.543	1	.000	.966
	employ	051	.010	25.767	1	.000	.950
	reside	103	.046	5.037	1	.025	.902
	equip	-1.948	.381	26.180	1	.000	.143
	callcard	.777	.151	26.451	1	.000	2.175
	longmon	233	.022	115.619	1	.000	.792
	equipmon	042	.011	15.377	1	.000	.959
	multline	.612	.145	17.854	1	.000	1.844
	voice	501	.157	10.197	1	.001	.606
	internet	362	.160	5.114	1	.024	.697
	callid	464	.148	9.790	1	.002	.629
	ebill	399	.156	6.557	1	.010	.671

Figure 348. Variables de l'équation (étape 12 uniquement)

Le modèle final contient les variables address, employ, reside, equip, callcard, longmon, equipmon, multline, voice, internet, callid, et ebill. Pour comprendre les effets des prédicteurs individuels, examinez Exp(B) qui peut être interprété comme la modification prédite du risque d'augmentation d'unités dans le prédicteur.

- La valeur de Exp(B) pour la variable *address* signifie que le risque d'attrition est réduit de 100% (100%0,966)=3,4% pour chaque année où le client a vécu à la même adresse. Le risque d'attrition pour un client ayant vécu à la même adresse pendant cinq ans est réduit de 100%–(100%×0.966⁵)=15.88%.
- La valeur de Exp(B) pour la variable *callcard* signifie que le risque d'attrition pour un client ne s'étant pas abonné au service de carte téléphonique est 2,175 fois plus élevé que celui pour un client s'étant abonné à ce service. Souvenez-vous que selon les codages de variable catégorielle *Non* = 1 pour la régression.
- La valeur de Exp(B) pour la variable *internet* signifie que le risque d'attrition pour un client ne s'étant pas abonné au service Internet est 0,697 fois plus élevé que celui pour un client s'étant abonné à ce service. Ce chiffre est quelque peu inquiétant car il indique que les clients abonnés au service quittent plus rapidement l'entreprise que ceux n'étant pas abonnés.

	30	Score	df	Siq.
Step 12	age	.122	1	.726
	marital	.648	1	.421
	income	1.476	1	.224
	ed	6.328	4	.176
	ed(1)	.007	1	.934
	ed(2)	.203	1	.652
	ed(3)	.835	1	.361
	ed(4)	5.773	1	.016
	retire	.013	1	.908
	gender	.214	1	.644
	tollfree	3.243	1	.072
	wireless	.668	1	.414
	tollmon	.000	1	.987
	cardmon	3.163	1	.075
	wiremon	1.084	1	.298
	pager	1.808	1	.179
	callwait	.266	1	.606
	forward	2.201	1	.138
	confer	2.568	1	.109
	custcat	.864	3	.834
	custcat(1)	.466	1	.495
	custcat(2)	.450	1	.502
	custcat(3)	.019	1	.889

Figure 349. Variables absentes du modèle (étape 12 uniquement)

Les variables absentes du modèle ont toutes des statistiques de score avec des valeurs de signification supérieures à 0,05. Cependant, les valeurs de signification de *tollfree* et de *cardmon*, bien qu'égales ou supérieures à 0,05, en sont relativement proches. Il serait intéressant d'étudier cette question plus avant.

Moyennes des covariables

	Mean		
age	41.684		
marital	.505		
address	11.551		
income	77.535		
ed(1)	.204		
ed(2)	.287		
ed(3)	.209		
ed(4)	.234		
employ	10.987		
retire	.953		
gender	.483		
reside	2.331		
tollfree	.526		
equip	.614		
callcard	.322		
wireless	.704		
longmon	11.723		
tollmon	13.274		
equipmon	14.220		
cardmon	13.781		
wiremon	11.584		
multline	.525		
voice	.696		
pager	.739		
internet	.632		
callid	.519		
callwait	.515		
forward	.507		
confer	.498		
ebill	.629		
custcat(1)	.266		
custcat(2)	.217		
custcat(3)	.281		

Figure 350. Moyennes des covariables

Ce tableau affiche la valeur moyenne de chaque variable de prédicteur. Ce tableau est une référence utile pour l'examen des tracés de survie qui sont construits pour les valeurs moyennes. Cependant, veuillez noter que le client "moyen" n'existe pas si l'on examine les moyennes des variables indicateur pour les prédicteurs indépendants. Même avec tous les prédicteurs d'échelle, il est peu probable que vous trouviez un client dont les valeurs de covariables soient toutes proches de la moyenne. Si vous souhaitez visualiser la courbe de survie d'une observation spécifique, vous pouvez modifier les valeurs de covariables qui déterminent le tracé de la courbe de survie dans la boîte de dialogue Tracés. Si vous souhaitez visualiser la courbe de survie d'une observation spécifique, vous pouvez modifier les valeurs de covariables qui déterminent le tracé de la courbe de survie dans la boîte de dialogue Tracés. Si vous souhaitez visualiser la courbe de survie d'une observation spécifique, vous pouvez modifier les valeurs de covariables qui déterminent le tracé de la courbe de survie dans le groupe Tracés de la boîte de dialogue Sorties avancées.

Courbe de survie



Figure 351. Courbe de survie pour le client "moyen"

La courbe de survie de base est un affichage visuel de la durée jusqu'à l'attrition pour le client "moyen" prédite par le modèle. L'axe horizontal indique la durée jusqu'à l'événement. L'axe vertical indique la probabilité de survie. Ainsi, tous les points de la courbe de survie indiquent la probabilité que le client "moyen" reste client cette durée passée. Après 55 mois, la courbe de survie devient plus irrégulière. Il existe moins de clients restés clients de l'entreprise pendant aussi longtemps et les informations disponibles sont plus rares, ce qui crée une courbe en dents de scie.

Courbe de risque



Figure 352. Courbe de risque pour le client "moyen"

La courbe de risque de base est un affichage visuel des risques cumulatifs d'attrition pour le client "moyen" prédits par le modèle. L'axe horizontal indique la durée jusqu'à l'événement. L'axe vertical indique les risques cumulatifs, égaux au log négatif de la probabilité de survie. Après 55 mois, la courbe de risque, comme la courbe de survie, devient plus irrégulière pour la même raison.

Evaluation

Les méthodes de sélection étape par étape garantissent que votre modèle n'ait que des prédicteurs "significatifs en termes de statistiques", mais ne garantissent pas que le modèle soit approprié pour prédire la cible. Pour ceci, vous devez analyser les enregistrements évalués.

😡 Evalua	ation					×
cox	ip File	N <u>G</u> enerate	e 🌔 Pre	view) 🔬	0	-0
Settings	Advanced	Summary	Annotations			
Predict sur	vival at futur	e times speci	ified as:			
O Regular	intervals	Time	e interval:		1.0 🔷	
		Num	iber of time pe	riods to score:	1 🜩	
Time fiel	d	🔗 te	enure			-1
Past surviv	al time:					-
V Append	all probabilit	ies				
Calculat	e cumulative	hazard func	tion			
ОК	Cancel					Reset

Figure 353. Nugget de Cox : Onglet Paramètres

- 1. Placez le nugget du modèle dans le canevas et liez-le au noeud source, ouvrez le nugget et cliquez sur l'onglet Paramètres.
- 2. Sélectionnez un **Champ temporel** et définissez la *durée d'affectation*. Chaque enregistrement sera évalué en fonction de sa durée d'affectation.
- 3. Sélectionnez Ajouter toutes les probabilités.

Cette action crée des évaluations qui utilisent 0,5 comme césure de l'attrition des clients ; si leur propension à quitter le service est supérieure à 0,5, ils sont évalués comme clients perdus. Ce chiffre n'est pas un chiffre magique et une césure différente peut offrir de meilleurs résultats. Utiliser le noeud Evaluation est une façon de choisir la césure.

🚱 [\$C-churn-1]	X								
	0								
Plot Options Appearance Output Annotations									
Chart type: 🔘 Gains 🔘 Response 🔍 Lift	🔘 Profit 🛛 🔍 R O I								
Cumulative plot Include baseline Include baseline	est line								
Find predicted/predictor fields using:	Models								
Model output field metadata									
Field name format (for example, '\$ <x>-<target field="">')</target></x>									
Other Score Fields									
Plot score fields									
Target	-1								
Separate by partition									
Plot: Percentiles 💌									
Style: Line Point 									
Costs: Fixed 5.0 No No No No No No No No No N	riable								
Revenue: 🔘 Fixed 10.0 🖨 🔘 Va	riable								
Weight:	riable								
OK Run Cancel	Apply Reset								

Figure 354. Noeud Evaluation : Onglet Tracé

- 4. Liez un noeud évaluation au nugget de modèle ; dans l'onglet Tracé, sélectionnez **Inclure la meilleure ligne**.
- 5. Cliquez sur l'onglet **Options**.

SC-chur	n-1]			D	<
				0	
Plot Options	Appearance Outpu	Annotations			
🔲 User define	d hit				
Condition:					
🔽 User define	d score				
Expression:	'\$CP-1-1')
nclude bus	ness rule				
Condition:					
Export resu	tts to file				
Filename:	output.txt				
Delimiter:					
Include field	names 🛛 🐨 New line a	after each recor	d		
					_
ОК	Run Cancel			Apply Reset	

Figure 355. Noeud Evaluation : Onglet Options

- 6. Sélectionnez **Score défini par l'utilisateur** et saisissez l'expression '\$CP-1-1'. C'est un champ généré par le modèle qui correspond à la propension à l'attrition.
- 7. Cliquez sur Exécuter.



Figure 356. Graphique de gain

Le graphique des gains cumulés montre le pourcentage du nombre total d'observations dans une modalité donnée obtenu en ciblant un pourcentage du nombre total d'observations. Par exemple, un point de la courbe est à (10%, 15%), ce qui signifie que si vous évaluez un jeu de données avec le modèle et triez toutes les observations en fonction de la propension à l'attrition prédite, les premiers 10% devraient contenir environ 15% de toutes les observations qui correspondent à la catégorie 1 (clients perdus). De la même façon, les premiers 60% contiennent environ 79,2% des clients perdus. Si vous sélectionnez 100% du jeu de données évalué, tous les clients perdus sont dans ce jeu de données.

La diagonale est la courbe "de référence" ; si vous sélectionnez au hasard 20% des enregistrements du jeu de données évalué, vous devriez "obtenir" environ 20% de tous les enregistrements qui correspondent à la catégorie 1. La "meilleure" ligne représente la courbe d'un modèle "parfait" qui attribue un score de propension à l'attrition plus élevé à tous les clients perdus plutôt qu'à tous les clients non perdus. Vous pouvez utiliser le graphique des gains cumulés pour sélectionner une césure de classement en choisissant un pourcentage correspondant à un gain souhaitable, puis en associant ce pourcentage à la valeur de césure appropriée.

Ce qui constitue un gain « souhaitable » dépend du coût des erreurs de type I et de type II. C'est-à-dire, combien coûte le classement d'un client perdu en client retenu (Type I) ? Combien coûte le classement d'un client retenu en client perdu (Type II) ? Si la rétention du client est la préoccupation principale, il vous faut alors diminuer votre erreur de Type I ; sur le graphique des gains cumulatifs, cela peut correspondre à une augmentation de l'assistance clientèle pour les clients faisant partie des premiers 60% de la propension prédite de 1, qui rassemblent 79,2% des clients perdus potentiels mais qui coûtent cher en temps et en ressources qui pourraient être utilisés pour acquérir de nouveaux clients. Si la baisse du coût du maintien de votre base de clientèle actuelle est votre priorité, diminuez alors votre erreur de Type II. Sur le graphique, cela peut correspondre à une augmentation de l'assistance clientèle pour les premiers 20% qui rassemblent 32,5% des clients perdus. Généralement, ces deux préoccupations sont importantes et il est nécessaire de choisir une règle de décision qui classifie les clients et qui offre le meilleur compromis entre la sensibilité et la spécificité.

Sort Sort	X
	0.0
Settings Optimization Annotations	
Sort by:	
Field	Order
\$CP-1-1	Descending
	<u>^</u>
	+
	+
Default sort order: O Ascending O Descending	
OK Cancel	Apply Reset

Figure 357. Noeud Trier : Onglet Paramètres

- 8. Imaginons que vous avez décidé que 45,6% est un gain acceptable ce qui correspond à utiliser les premiers 30% des enregistrements. Pour trouver une césure de classification appropriée, liez un noeud Trier au nugget de modèle.
- 9. Dans l'onglet Paramètres, choisissez de trier par *\$CP-1-1* dans l'ordre décroissant puis cliquez sur **OK**.

違 <u>F</u> ile		Edit 🖔 🧐	Generate		@ >
Table	Annota	ations			
	Irn \$	6C-churn-1	\$CP-churn-1	\$CP-0-1	\$CP-1-1
292	0	1	0.744	0.744	0.256
293	0	í	0.745	0.745	0.255
294	0		0.745	0.745	0.255
295	0		0.746	0.746	0.254
296	0		0.748	0.748	0.252
297	0		0.749	0.749	0.251
298	0		0.749	0.749	0.251
299	0		0.750	0.750	0.250
300	0	I	0.752	0.752	0.248
301	0		0.752	0.752	0.248
302	0		0.754	0.754	0.246
303	0	I	0.754	0.754	0.246
304	0		0.755	0.755	0.245
305	0	i	0.756	0.756	0.244
306	0		0.757	0.757	0.243
307	0		0.757	0.757	0.243
308	0	i	0.758	0.758	0.242
309	0	1	0.759	0.759	0.241
310	0	1	0.761	0.761	0.239
311	0	1	0.762	0.762	0.238
	4				

Figure 358. Table

- 10. Reliez un noeud Table au noeud Trier.
- 11. Ouvrez le noeud Table, puis cliquez sur Exécuter.

En faisant défiler les sorties vers le bas, vous pouvez voir que la valeur de *\$CP-1-1* est de 0,248 pour le 300ème enregistrement. L'utilisation de 0.248 comme césure de classification devrait résulter approximativement en 30% des clients évalués comme clients perdus, capturant à peu près 45% du nombre total des clients perdus.

Suivi du nombre prévu de clients retenus

Lorsque le modèle vous convient, effectuez un suivi du nombre prévu de clients du jeu de données qui sont retenus pendant les deux prochaines années. Les valeurs nulles, qui correspondent à des clients dont la durée d'affectation totale (temps futur + *durée d'affectation*) est en-dehors de la plage des durées de survie pour les données utilisées pour former le modèle, représentent un défi intéressant. Une façon de traiter ces valeurs est de créer deux ensembles de prédictions, un dans lequel on considère que les valeurs nulles ont été perdues et l'autre dans lequel on considère que ces valeurs ont été retenues. Ainsi, vous pouvez établir les limites supérieures et inférieures du nombre prévu de clients retenus.

😡 Predie	ction					
cox	違 <u>F</u> ile	🏷 Generat	te 🌔 Prev	view) 📳	0	
Settings	Advanced	Summary	Annotations			
Predict sur	vival at futur	e times spec	cified as:			
🔘 Regular	intervals	Tim	e interval:		1.0 🚔	
		Nun	nber of time pe	riods to score:	24 荣	
🔘 Time fie	ld					-1
Past surviv	al time:	🥜 te	enure			
📝 Append	l all probabiliti	es				
🗾 Calculat	te cumulative	hazard fund	ction			
ОК	Cancel				Apply	Reset

Figure 359. Nugget de Cox : Onglet Paramètres

- 1. Double-cliquez sur le nugget de modèle dans la palette Modèles (ou copiez-collez le nugget sur l'espace de travail du flux) et joignez le nouveau nugget au noeud Source.
- 2. Ouvrez le nugget dans l'onglet Paramètres.
- **3**. Vérifiez que **Intervalles réguliers** est sélectionné et spécifiez 1,0 comme intervalle de temps et 24 comme nombre de périodes à évaluer. Chaque enregistrement sera ainsi évalué pendant chacun des 24 mois suivants.
- 4. Sélectionnez le champ *durée d'affectation* pour spécifier la durée de survie passée. L'algorithme d'évaluation prendra en compte la durée de survie de chaque client comme client de l'entreprise.
- 5. Sélectionnez Ajouter toutes les probabilités.

🖸 Lower Estimate 🛛 🔯								
>	Preview)				0		
Settings	Annotations							
Key fields:						📃 Keys ar	e contiguous	
Aggregate f	ields:						—	
Field		Sum	Mean	Min	Max	SDev		
\$CP-0-1		-						
\$CP-0-10		-						
\$CP-0-11		-						
\$CP-0-12		-						
\$CP-0-13		\checkmark						
\$CP-0-14		-					-	
Default mode: Sum Mean Min Max SDev								
New field na	ame extension	r		Add as:	🔘 Suffix	O Prefix		
🔲 Include r	ecord count i	n field Rec	ord_Count					
ОК	OK Cancel Apply Reset							

Figure 360. Noeud agrégé : Onglet Paramètres

- 6. Liez un noeud agrégé au nugget de modèle ; dans l'onglet Paramètres, désélectionnez le mode par défaut **Moyenne**.
- 7. Sélectionnez \$*CP-0-1* à \$*CP-0-24*, les champs de forme \$*CP-0-n*, étant les champs à agréger. Pour faciliter cette action, triez les champs par nom (c'est-à-dire par ordre alphabétique) dans la boîte de dialogue Sélectionner les champs.
- 8. Désélectionnez Inclure le comptage des enregistrements dans le champ.
- 9. Cliquez sur OK. Ce noeud crée les prévisions de "limite inférieure".

💟 Nulls Stay	X
Settings Annotations	0
Fill in fields:	
Condition:	
@BLANK(@FIELD)	
Replace with:	
1	
OK Cancel	Apply Reset

Figure 361. Noeud remplissage : Onglet Paramètres

- 10. Liez un noeud remplissage au nugget Coxreg auquel le noeud agrégé vient d'être lié ; dans l'onglet Paramètres, sélectionnez \$*CP-0-1* à \$*CP-0-24*, les champs de forme \$*CP-0-n*, comme champs à remplir. Pour faciliter cette action, triez les champs par nom (c'est-à-dire par ordre alphabétique) dans la boîte de dialogue Sélectionner les champs.
- 11. Choisissez de remplacer les Valeurs nulles par la valeur 1.
- 12. Cliquez sur OK.

🖓 Upper Estimate 🛛 🛛 🔯							
*	Preview					0	
Settings	Annotations						
Key fields:						📕 Keys ar	e contiguous
Aggregate	fields:						×
Field		Sum	Mean	Min	Max	SDev	
\$CP-0-1		-					
\$CP-0-10		-					
\$CP-0-11		-					
\$CP-0-12		-					
\$CP-0-13		-					
\$CP-0-14		-					*
Default mode: Sum Mean Min Max SDev							
	record count i	n field Rec	ord_Count	Hun do.	3 Gamix	S HOIX	
ок	Cancel					Appl	y <u>R</u> eset

Figure 362. Noeud agrégé : Onglet Paramètres

- **13**. Liez un noeud agrégé au noeud remplissage ; dans l'onglet Paramètres, désélectionnez le mode par défaut **Moyenne**.
- 14. Sélectionnez *\$CP-0-1* à *\$CP-0-24*, les champs de forme *\$CP-0-n*, étant les champs à agréger. Pour faciliter cette action, triez les champs par nom (c'est-à-dire par ordre alphabétique) dans la boîte de dialogue Sélectionner les champs.
- 15. Désélectionnez Inclure le comptage des enregistrements dans le champ.
- 16. Cliquez sur OK. Ce noeud crée les prévisions de "limite supérieure".
| Months | | X |
|---|-----------------|---------------------------------------|
| Rreview | | |
| Filter Annotations | | |
| 7 • • | Fields: | 24 in, 0 filtered, 24 renamed, 24 out |
| Field - | Filter | Field |
| \$CP-0-1_Sum | \rightarrow | 1 🖊 |
| \$CP-0-2_Sum | \rightarrow | 2 |
| \$CP-0-3_Sum | \rightarrow | 3 |
| \$CP-0-4_Sum | \rightarrow | 4 |
| \$CP-0-5_Sum | \rightarrow | 5 |
| \$CP-0-6_Sum | \rightarrow | 6 |
| \$CP-0-7_Sum | \rightarrow | 7 |
| \$CP-0-8_Sum | \rightarrow | 8 |
| \$CP-0-9_Sum | \rightarrow | 9 |
| \$CP-0-10_Sum | \rightarrow | 10 |
| View current fields View Cancel | unused field se | attings |

Figure 363. Noeud Filtrer : Onglet Paramètres

- 17. Reliez un noeud Ajouter aux deux noeuds Agrégé puis reliez un noeud filtre au noeud Ajouter.
- **18**. Dans l'onglet Paramètres du noeud Filtrer, renommez les champs de *1* à 24. En utilisant un noeud Transposer, les noms de ces champs deviendront des valeurs de l'axe *x* dans les graphiques en aval.

🙀 Transpose			
Preview)	0	
Settings Annotations			
New field names:			
Our Dese prefix	Field	Number of new fields:	2 🖨
Read from field			-
Read Values	New Field Names		
Transpose: 🔘 All nur	Maximum number of values to read: neric © All string © Custom	500	
Fields:		×	
Row ID name: ID			
OK Cancel		App	oly <u>R</u> eset

Figure 364. Noeud Transposer : Onglet Paramètres

- 19. Reliez un noeud Transposer au noeud Filtrer.
- 20. Saisissez 2 comme nombre des nouveaux champs.

😡 Labels		
Review]		0
Filter Annotations		
7- 📭	Field	s: 3 in, 0 filtered, 3 renamed, 3 out
Field -	Filter	Field
ID	\rightarrow	Months
Field1	\rightarrow	Lower Estimate
Field2	\rightarrow	Upper Estimate
View current fields O View	v unused field sett	ngs
OK Cancel		Apply Reset

Figure 365. Noeud Filtrer : Onglet Filtrer

- 21. Reliez un noeud Filtrer au noeud Transposer.
- **22**. Dans l'onglet Paramètres du noeud Filtrer, changez le nom de *ID* en *Mois*, *Champ1* en *Estimation inférieure*, et *Champ2* en *Estimation supérieure*.

🚱 Estimated Numbers	
Plot Appearance Output Annotations	0
X field	
Y fields:	
Overlay	
Panet: Animation:	J
Normalize	
Overlay function y = :	
When number of records greater than: 2000 🗬	
⊚ Bin ◯ Sample ◯ Use all data	
OK 🕨 Run Cancel	Apply Reset

Figure 366. Noeud Courbes : Onglet Tracé

- 23. Reliez un noeud Courbes au noeud Filtrer.
- 24. Dans l'onglet Tracé, sélectionnez *Months* (Mois) comme champ X, et *Lower Estimate* (Estimation inférieure) et *Upper Estimate* (Estimation supérieure) comme champs Y.

Estimated Numbers	
Plot Appearance Output Annotations	
Title: Number of Customers	
Subtitle:	
Caption: Estimates the number of customers retained	
X label: 🔘 Auto 🔘 Custom	
Y label: 🔘 Auto 🔘 Custom	
☑ Display gridline	
OK 🕨 Run Cancel	Apply Reset

Figure 367. Noeud Courbes : Onglet Apparence

- 25. Cliquez sur l'onglet Apparence.
- 26. Saisissez Nombre de clients comme titre.
- 27. Saisissez Estimations du nombre de clients retenus comme légende.
- 28. Cliquez sur Exécuter.



Estimates the number of customers retained

Figure 368. Courbes estimant le nombre de clients retenus

Les limites supérieures et inférieures du nombre estimé de clients retenus sont représentées. La différence entre les deux lignes est le nombre de clients évalués comme nuls et, par conséquent, dont le statut est fortement incertain. Au fil du temps, le nombre de ces clients augmente. Après 12 mois, vous devriez retenir entre 601 et 735 des clients d'origine du jeu de données ; après 24 mois, entre 288 et 597.

💟 Unknown %	×
	0
Derive as: Formula	
Settings Annotations	
Mode: 💿 Single 🔘 Multiple	
Derive field:	
Unknown %	
Derive as: Formula T Field type: Continuous T	
(100 * ('Upper Estimate' - 'Lower Estimate')) / 'Lower Estimate'	
OK Cancel	Apply Reset

Figure 369. Noeud Dériver : Onglet Paramètres

- **29**. Pour examiner de nouveau combien les estimations du nombre de clients retenus sont incertaines, reliez un noeud Dériver au noeud Filtrer.
- 30. Dans l'onglet Paramètres du noeud Dériver, saisissez % inconnu comme champ de calcul.
- 31. Sélectionnez Continu comme type de champ.
- 32. Saisissez la formule (100 * ('Estimation supérieure' 'Estimation inférieure')) / 'Estimation inférieure'. % *inconnu* est le nombre de clients "dans le doute" sous la forme d'un pourcentage de l'estimation inférieure.
- 33. Cliquez sur OK.

🖸 Unpredicted Cases 🛛 🛛 🛛
X: Months Y: Unknown %
Plot Options Appearance Output Annotations
L X field: Months Y field: V field:
Overlay
Color: Size: Shape:
Panel: Animation: Transparency:
Overlay type: None
O Smoother
© Function y =
OK Run Cancel Apply Reset

Figure 370. Noeud Tracé : Onglet Tracé

- 34. Reliez un noeud Nuage au noeud Dériver.
- **35**. Dans l'onglet Tracé du noeud Tracé, sélectionnez *Months* (Mois) comme champ X et *Unknown* % (% inconnu) comme champ Y.
- 36. Cliquez sur l'onglet Apparence.

X: Months Y: Unknown % Plot Options Appearance Output Annotations Title: Unpredictable Customers as % of Predictable Customers Subtitle:	💟 Unpred	licted Cases					X
X: Months Y: Unknown % Plot Options Appearance Output Annotations Title: Unpredictable Customers as % of Predictable Customers Image: Caption: Image: Caption: Image: Custom X label: Image: Auto Image: Custom Image: Custom Image: Custom Image: Custom Y label: Image: Auto Image: Custom Image: Custom Image: Custom Image: Custom						0	
Plot Options Appearance Output Annotations Title: Unpredictable Customers as % of Predictable Customers Subtitle:		: Months	Y: Uni	known %			
Title: Unpredictable Customers as % of Predictable Customers Subtitle:	Plot Optio	Appearance	Output	Annotations			
Subtitle:	Title: Unj	predictable Custome	ers as %	of Predictable	Customers		
Caption:X label: Auto Custom Y label: Auto Custom	Subtitle:						
X label: Auto Custom Y label: Auto Custom 	Caption:						
Y label: 🔘 Auto 🔘 Custom	X label:	🖲 Auto 🔘 Custo	m				
	Y label:	🔘 Auto 🔘 Custo	m				
Z label: 🔘 Auto 🔘 Custom	Z label:	Auto O Custo	m				
Display gridline	Display	gridline					

Figure 371. Noeud Tracé : Onglet Apparence

- 37. Saisissez le titre Clients imprévisibles comme % des clients prévisibles.
- **38**. Exécutez le noeud.



Figure 372. Tracé des clients imprévisibles

Pendant la première année, le pourcentage des clients imprévisibles augmente assez régulièrement, mais explose pendant la deuxième année jusqu'à ce que, le 23ème mois, le nombre de clients avec des valeurs nulles dépasse le nombre prévu de clients retenus.

Scoring

Lorsque votre modèle vous convient, évaluez les clients afin d'identifier les individus les plus susceptibles d'attrition l'année suivante, par trimestre.



Figure 373. Nugget Coxreg : Onglet Paramètres

- 1. Joignez un troisième nugget de modèle au noeud Source et ouvrez le nugget de modèle.
- **2**. Vérifiez que **Intervalles réguliers** est sélectionné et spécifiez 3.0 comme intervalle de temps et 4 comme nombre de périodes à évaluer. Chaque enregistrement sera ainsi évalué pendant les quatre trimestres suivants.
- **3**. Sélectionnez le champ *durée d'affectation* pour spécifier la durée de survie passée. L'algorithme d'évaluation prendra en compte la durée de survie de chaque client comme client de l'entreprise.
- 4. Sélectionnez **Ajouter toutes les probabilités**. Ces champs supplémentaires faciliteront le tri des enregistrements et leur affichage dans un tableau.

Select		
-7>	Preview	0 - 🗖
X		
Settings 4	nnotations	
Mode:	linclude 🔘 Discard	
	churn = 0	
Condition:		
окс	ancel	Apply

Figure 374. Noeud Sélectionner : Onglet Paramètres

5. Reliez un noeud Sélectionner au nugget de modèle ; dans l'onglet Paramètres, saisissez la condition churn=0. Cette action supprime les clients déjà perdus du tableau de résultats.

Churn					X
Preview)				0	
Derive as: Con	ditional				
Settings Annotations					
	Mode:	🔘 Single 🧕	Multiple		
Derive from:					
♦ \$CP-1-1					
\$CP-1-2 \$\$CP-1-3					, ×
Field name extension:	churn		Add as:	Suffix	O Prefix
Derive as: Conditional 🔽		TIP: R	efer to selecte	d fields by us	sing @FIELD
Field type: 🔓 Flag	*				
lf:					
@FIELD>0.248					
Theo:					
1					
Else:					
OK Cancel				(lqq <u>A</u>	y <u>R</u> eset

Figure 375. Noeud Dériver : Onglet Paramètres

- 6. Reliez un noeud Dériver au noeud Sélectionner ; dans l'onglet Paramètres, sélectionnez le mode **Multiple**.
- 7. Choisissez de calculer les champs \$*CP-1-1* à \$*CP-1-4*, les champs de forme \$*CP-1-n*, et saisissez le suffixe _attrition à ajouter. Pour faciliter cette action, triez les champs par nom (c'est-à-dire par ordre alphabétique) dans la boîte de dialogue Sélectionner les champs.
- 8. Choisissez de calculer le champ comme champ Conditionnel.
- 9. Sélectionnez Indicateur comme niveau de mesure.
- 10. Saisissez @FIELD>0.248 comme condition If (Si). Veuillez noter qu'il s'agit là de la césure de classification identifiée pendant l'évaluation.
- 11. Saisissez 1 comme expression Then (Alors).
- 12. Saisissez 0 comme expression Else (Sinon).
- 13. Cliquez sur OK.

🚰 Sort		
Preview	0	
Settings Optimization Annotations		
Field	Order	
\$CP-1-1_churn	🔻 Descending 🧧	
\$CP-1-2_churn	Descending	×
\$CP-1-3_churn	Descending	
\$CP-1-4_churn	Descending	
\$CP-1-1	Descending	÷
\$CP-1-2	Descending	
\$CP-1-3	🔻 Descending	-
Default sort order: O Ascending O Descending		Reset

Figure 376. Noeud Trier : Onglet Paramètres

14. Reliez un noeud Trier au noeud Dériver ; dans l'onglet Paramètres, choisissez de trier par \$*CP-1-1_attrition* à \$*CP-1-4-attrition* puis par \$*CP-1-1* à \$*CP-1-4*, dans l'ordre décroissant. Les clients dont l'attrition a été prévue figureront en premier.

🙀 Fiel	d Reorder		
Reorde	Appdations		0
Cust	om Order	O Automatic Sort	
Type:	Name: 🔊 Si		
Туре	Field	Storage	
8	\$CP-1-1 churn	? (Unknown)	
	\$CP-1-1	Real	
8	\$CP-1-2_churn	(Unknown)	
	\$CP-1-2	🛞 Real	^
8	\$CP-1-3_churn	💈 (Unknown)	
	\$CP-1-3	🛞 Real	
8	\$CP-1-4_churn	🦚 (Unknown)	→
	\$CP-1-4	🛞 Real	
Clear I Note: Fi	Jnused	e are not reordered.	
ок	Cancel		Apply Reset

Figure 377. Noeud Réorganiser : onglet Réorganiser

15. Reliez un noeud Réorganiser au noeud Trier ; dans l'onglet Réorganiser, choisissez de placer les champs *\$CP-1-1_attrition* à *\$CP-1-4* devant les autres champs. Ceci est une option qui permet simplement de faciliter la lecture du tableau de résultats. Utilisez les boutons pour déplacer les champs comme indiqué dans le schéma.

🛄 Table	e (50 fiel	ds, 726 record	ls)						
😂 File	📄 Edit	🏷 <u>G</u> enerate		22					0>
Table ,	Annotations				201				
	\$CP-1-1_	churn \$CP-1-1	\$CP-1-2_churn	\$CP-1-2	\$CP-1-3_churn	\$CP-1-3	\$CP-1-4_churn	\$CP-1-4	tenur
255	0	0.032	0	0.075	0	0.147	1	0.298	49 4
256	0	0.027	0	0.064	0	0.127	1	0.260	49
257	0	0.023	0	0.130	0	0.233	1	0.308	53
258	0	0.021	0	0.127	0	0.239	1	0.320	54
259	0	0.021	0	0.125	0	0.237	1	0.318	54
260	0	0.021	0	0.053	0	0.198	1	0.331	50
261	0	0.021	0	0.053	0	0.196	1	0.329	50
262	0	0.020	0	0.050	0	0.189	1	0.317	50
263	0	0.017	0	0.043	0	0.163	1	0.278	50
264	0	0.015	0	0.039	0	0.148	1	0.253	50
265	0	0.197	0	0.197	0	\$null\$	0	\$null\$	66
266	0	0.109	0	0.109	0	\$null\$	0	\$null\$	66
267	0	0.101	0	0.214	0	\$null\$	0	\$null\$	65
268	0	0.081	0	0.137	0	0.194	0	0.245	23
269	0	0.074	0	0.159	0	\$null\$	0	\$null\$	65
270	0	0.070	0	0.116	0	0.158	0	0.237	28
271	0	0.070	0	0.128	0	0.189	0	0.234	45
272	0	0.062	0	0.105	0	0.151	0	0.191	23
273	0	0.062	0	0.130	0	0.163	0	0.212	44
274	0	0.061	0	0.123	0	0.182	0	0.241	4
					10 km				1
									U

Figure 378. Tableau présentant les scores des clients

16. Reliez un noeud Table au noeud Réorganiser et exécutez-le.

L'attrition de 264 clients est prévue pour la fin de l'année, 184 à la fin du troisième trimestre, 103 à la fin du deuxième et 31 à la fin du premier. Veuillez noter que sur deux clients, celui avec la plus forte propension à l'attrition pendant le premier trimestre n'a pas nécessairement la plus forte propension à l'attrition pendant les trimestres suivants ; examinez par exemple les enregistrements 256 et 260. Ceci est probablement dû à la forme de la fonction des risques des mois suivant la durée d'affectation actuelle du client ; par exemple, les clients qui ont rejoint l'entreprise en raison d'une promotion sont plus susceptibles de partir plus rapidement que les clients ayant rejoint l'entreprise sur recommandation personnelle, mais s'ils ne partent pas, ils peuvent alors être plus fidèles pour leur durée d'affectation restante. Réorganiser de nouveau les clients pour obtenir des vues différentes des clients les plus susceptibles de quitter.

🔟 Table	e (50 fields, 7	26 record	ds)						- 0
🐞 <u>F</u> ile	📄 Edit 🛛 🐑	Generate		22					@ >
Table ,	Annotations								
	\$CP-1-1_churn	\$CP-1-1	\$CP-1-2_churn	\$CP-1-2	\$CP-1-3_churn	\$CP-1-3	\$CP-1-4_churn	\$CP-1-4	tenur
707	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	71 4
708	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	71
709	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	71
710	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	72
711	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	71
712	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	72
713	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	72
714	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	72
715	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	70
716	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	70
717	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	71
718	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	72
719	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	72
720	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	72
721	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	72
722	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	71
723	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	70
724	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	71
725	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	70
726	0	\$null\$	0	\$null\$	0	\$null\$	0	\$null\$	72
for a second second									•
									OF

Figure 379. Tableau présentant les clients avec des valeurs nulles

Les clients avec des valeurs nulles prédites se trouvent au bas de la table. Il s'agit de clients dont la durée d'affectation totale (temps futur + *durée d'affectation*) se trouve en-dehors de la plage des durées de survie des données utilisées pour former le modèle.

Récapitulatif

L'utilisation de la régression de Cox vous a permis de trouver un modèle approprié pour la durée jusqu'à l'attrition, de représenter le nombre prévu de clients retenus pendant les deux prochaines années et d'identifier les clients individuels les plus susceptibles de quitter au cours de l'année suivante. Remarque : même si ce modèle semble acceptable, il n'est peut-être pas le meilleur modèle. Idéalement, vous devriez comparer ce modèle, obtenu à partir de la méthode Pas à pas ascendante, à un modèle qui utilise la méthode Pas à pas descendante.

Vous trouverez des explications sur le fondement mathématique des méthodes de modélisation utilisées dans IBM SPSS Modeler dans le *guide des algorithmes IBM SPSS Modeler*.

Chapitre 27. Analyse d'un panier de courses (Induction de règle/C5.0)

Cet exemple se base sur des données fictives décrivant le contenu d'un panier à provisions (c'est-à-dire, un ensemble d'articles achetés en même temps) et sur les données personnelles de l'acheteur, lesquelles peuvent être collectées via un programme de fidélité. L'objectif est d'identifier des ensembles de consommateurs effectuant des achats similaires et pouvant être regroupés selon des caractéristiques démographiques telles que l'âge, les revenus, etc.

Cet exemple illustre deux phases du processus d'exploration de données :

- La modélisation de règles d'association et l'affichage des relations mettent en évidence les liens entre les articles achetés.
- L'induction d'une règle C5.0 permet d'établir un portrait des acheteurs des groupes de produits identifiés.

Remarque : Cette application ne fait pas directement appel à la modélisation prédictive, c'est pourquoi l'exactitude des modèles générés n'est pas mesurée et la distinction apprentissage/test n'est pas effectuée au cours du processus d'exploration de données.

Cet exemple utilise le flux nommé *baskrule* qui fait référence au fichier de données *BASKETS1n*. Ces fichiers sont disponibles dans le répertoire *Demos* de n'importe quelle installation d'IBM SPSS Modeler. Ce répertoire est accessible à partir du groupe de programmes IBM SPSS Modeler dans le menu Démarrer de Windows. Le fichier *baskrule* se trouve dans le répertoire des *flux*.

Accès aux données

Connectez un noeud Délimité au jeu de données *Panier* et choisissez de lire le nom des champs à partir du fichier. Connectez un noeud type à la source de données, puis connectez-le à un noeud table. Définissez le niveau de mesure du champ *No carte* sur *Sans type* (un numéro de carte de fidélité n'apparaissant qu'une fois dans le jeu de données, son utilisation ne présente pas d'intérêt particulier pour la modélisation). Sélectionnez *Nominal* comme niveau de mesure du champ *sexe* (afin que l'algorithme de modélisation Apriori ne considère pas le champ *sexe* comme un indicateur).



Figure 380. Flux baskrule (panier)

Exécutez le flux pour instancier le noeud type et affichez le tableau. Le jeu de données contient 18 champs, chacun représentant un panier.

Les 18 champs sont présentés sous les en-têtes suivants.

Récapitulatif du panier :

- No carte. Numéro de carte de fidélité de l'acheteur.
- montant. Prix total des articles du panier.
- paiement. Méthode de paiement.

Informations personnelles sur le détenteur de la carte :

- sexe
- locataire. Indique si le détenteur de la carte est locataire.
- revenu
- âge

Catégories de produits contenues dans le panier :

- fruits & légumes
- boucherie
- produits laitiers
- conserves légumes
- conserves viande
- surgelés
- bière
- vin
- boissons sucrées
- poisson
- confiseries

Identification des analogies entre les articles du panier

Vous devez tout d'abord obtenir un aperçu des analogies (associations) existant entre les articles du panier. Pour ce faire, utilisez un algorithme Apriori afin de produire des règles d'association. Sélectionnez les champs à utiliser au cours du processus de modélisation en définissant, dans le noeud type, le rôle de toutes les catégories de produits sur *Les deux* et toutes les autres rôles sur *Aucun*. (*Les deux* signifie que le champ peut figurer en tant qu'entrée ou en tant que sortie du modèle résultant.)

Remarque : Vous pouvez définir les options de plusieurs champs. Pour ce faire, maintenez la touche Maj enfoncée tout en sélectionnant les champs, puis spécifiez une option à partir des colonnes.

	eview				0
Types Format	Annotations	alues Clea	r Values	Clear All V	'alues
Field 🥅	Measurement	Values	Missing	Check	Role
A sex	💑 Nominal	F,M		None	○ None
homeown	🖁 Flag	YES/NO		None	S None
income	Continuous	[10200,3		None	○ None
age	Continuous	[16,50]		None	○ None
fruitveg	🖁 Flag	TÆ		None	🔘 Both 💌
freshmeat	Flag	TÆ		None	Input
dairy	🖁 Flag	TÆ		None	() Target
annedveg	🖁 Flag	TÆ		None	Roth
	U Eloa	TÆ		Nono	Bolli
View current	fields 🔘 View unu	sed field settin	gs		Partition

Figure 381. Sélection des champs pour la modélisation

Une fois les champs destinés à la modélisation spécifiés, connectez le noeud Apriori au noeud type, modifiez-le, sélectionnez l'option **Uniquement valeurs vraies pour indicateurs**, puis exécutez le noeud Apriori. Un modèle apparaît sur l'onglet Modèles situé dans la partie supérieure droite de la fenêtre des gestionnaires. Il contient des règles d'association que vous pouvez consulter en utilisant le menu contextuel et l'option **Parcourir**.

😭 11 fields) <u>G</u> enerate	ew) 🚯	
Model Settings Summ	nfidence %	■ - 7 🏦	3 of 3
Consequent	Antecedent	Support %	Confidence %
frozenmeal	beer cannedveg	16.7	87.425
cannedveg	beer frozenmeal	17.0	85.882
beer	frozenmeal cannedveg	17.3	84.393
OK Cancel			Apply Reset

Figure 382. Règles d'association

Ces règles montrent différentes associations entre les produits surgelés, les légumes en conserve et la bière. La présence de règles d'association d'ordre 2 du type :

surgelés -> bière bière -> surgelés

indique que l'affichage des relations (qui ne présente que ce type d'associations) pourrait permettre de dégager certaines tendances à partir de ces données.

Connectez un noeud relations au noeud type éditez le noeud relations, sélectionnez tous les champs correspondant au contenu du panier, cochez la case **Afficher uniquement les indicateurs ayant une valeur vraie** et exécutez le noeud relations.



Figure 383. Affichage des relations des associations entre les produits

La plupart des combinaisons de catégories de produits étant présentes dans plusieurs paniers, les liens forts figurant dans cette relation sont trop nombreux pour permettre de repérer les groupes d'acheteurs identifiés par le modèle.



Figure 384. Affichage des relations restreint

1. Pour spécifier les connexions faibles et les connexions fortes, cliquez sur la double flèche jaune de la barre d'outils. Ce bouton permet d'agrandir la boîte de dialogue, et d'afficher le récapitulatif et les commandes de sortie de la relation.

- 2. Sélectionnez La valeur de la taille est fort/normal/faible.
- 3. Définissez les liens faibles au-dessous de 90.
- 4. Définissez les liens forts au-dessus de 100.

Cet affichage généré met en évidence les groupes d'acheteurs suivants :

- Ceux qui consomment du poisson et des fruits et légumes (le groupe "santé" dans notre exemple).
- Ceux qui achètent du vin et des confiseries
- Ceux qui achètent de la bière, des plats surgelés et des légumes en conserve ("bière, petits pois et pizza")

Portrait des groupes d'acheteurs

Vous avez maintenant mis en évidence trois types de consommateur, regroupés en fonction des produits qu'ils achètent. Vous allez à présent les identifier plus en détail, en établissant leur profil démographique. Pour ce faire, vous pouvez associer chacun d'eux à un indicateur correspondant au groupe auquel il appartient, puis utiliser une règle C5.0 pour définir le profil de ces indicateurs.

Vous devez tout d'abord calculer un indicateur pour chaque groupe. Il peut être généré automatiquement en utilisant l'affichage des relations que vous venez de créer. A l'aide du bouton droit de la souris, cliquez sur le lien qui relie *fruits & légumes* et *poisson* pour le mettre en évidence; puis cliquez avec le bouton droit de la souris et sélectionnez **Générer le noeud dériver pour le lien**.



Figure 385. Calcul d'un indicateur pour chaque groupe d'acheteurs

Dans le noeud dériver généré, éditez le nom du champ et choisissez *santé*. Répétez l'opération avec le lien reliant *vin* à *confiseries*, et nommez le champ Calculer résultant *vin_choco*.

Pour le troisième groupe (impliquant trois liens), assurez-vous d'abord qu'aucun lien n'est sélectionné. Cliquez ensuite sur tous les liens du triangle *conserves légumes, bière* et *surgelés* pour les sélectionner tout en maintenant la touche MAJ enfoncée. (Assurez-vous d'être en mode interactif et non en mode d'édition.) Ensuite, choisissez les options de menu suivantes à partir de l'affichage des relations :

Générer > **Noeud dériver** (Et)

Modifiez le nom du champ Dériver résultant en bière_pizza_petits_pois.

Pour établir le profil de vos groupes d'acheteurs, connectez le noeud type existant à ces trois noeuds Dériver en série, puis connectez un autre noeud type. Dans le nouveau noeud type définissez le rôle *Aucun* pour tous les champs, sauf pour *montant*, *paiement*, *sexe*, *locataire*, *revenu* et âge, qui doivent être définis sur la valeur *Entrée*, et le groupe d'acheteurs défini (par exemple, *bière_pizza_petits_pois*), qui doit être défini sur la valeur *Cible*. Connectez un noeud C5.0, définissez le type de sortie sur **Ensemble de règles**, puis exécutez le noeud. Le modèle généré (pour *bière_pizza_petits_pois*) contient un profil démographique clair pour ce groupe d'acheteurs :

```
Rule 1 for T:
if sex = M
and income <= 16,900
then T
```

Pour appliquer la même méthode aux autres indicateurs de groupes d'acheteurs, sélectionnez-les comme sortie pour le second noeud type. D'autres profils peuvent être générés si vous utilisez Apriori au lieu de C5.0 dans ce contexte. Apriori permet également d'établir simultanément le profil de tous les indicateurs du groupe de clients, car il n'est pas limité à un seul champ de sortie.

Récapitulatif

Cet exemple illustre la manière dont IBM SPSS Modeler peut être utilisé pour mettre en évidence des associations, ou des liens, entre les éléments d'une base de données, par modélisation (à l'aide d'Apriori) et par visualisation (à l'aide de l'affichage des relations). Ces liens correspondent à des regroupements d'observations effectuées dans les données ; les groupes identifiés peuvent être analysés en détail et leur profil peut être établi par modélisation (à l'aide d'ensembles de règles C5.0).

Dans le domaine de la vente au détail, ces groupes peuvent permettre, par exemple, de cibler des offres spéciales afin d'obtenir de meilleurs résultats en termes de réponse au publipostage direct, ou de personnaliser la gamme de produits stockés par un magasin afin de répondre aux besoins de la clientèle, identifiée en fonction de caractéristiques démographiques.

Chapitre 28. Estimation des offres de nouveaux véhicules (KNN)

L'analyse du voisin le plus proche est une méthode de classification des observations en fonction de leur similarité avec les autres observations. En apprentissage automatique, elle a été développée comme une façon de reconnaître les configurations de données sans avoir à recourir à une correspondance exacte avec d'autres configurations ou observations stockées. Les observations semblables sont proches l'une de l'autre et les observations dissemblables sont éloignées l'une de l'autre. Par conséquent, la distance entre deux observations est une mesure de leur dissimilarité.

Les observations proches les unes des autres sont appelées « voisins ». Lorsqu'une nouvelle observation est présentée (traitée), sa distance de chacune des observations du modèle est calculée. Les classifications des observations les plus similaires « les plus proches voisins » sont mesurées et la nouvelle observation est placée dans la catégorie qui contient le plus grand nombre de voisins les plus proches.

Vous pouvez spécifier le nombre de voisins les plus proches à examiner, cette valeur est appelée k. Les illustrations montrent comment une nouvelle observation serait classée à l'aide de deux valeurs différentes de k. Lorsque k = 5, la nouvelle observation est placée dans la catégorie 1 car une majorité de voisins les plus proches appartiennent à la catégorie 1. Cependant, lorsque k = 9, la nouvelle observation est placée dans la catégorie 0 car une majorité de voisins les plus proches appartiennent à la catégorie 1.

L'analyse du voisin le plus proche peut également être utilisée pour calculer des valeurs pour une cible continue. Dans cette situation, la valeur cible de la médiane ou de la moyenne des voisins les plus proches est utilisée pour obtenir la valeur prédite de la nouvelle observation.

Un constructeur automobile a développé des prototypes pour deux nouveaux véhicules, une voiture et un camion. Avant de présenter les nouveaux modèles dans sa gamme, le constructeur souhaite déterminer les véhicules existant sur le marché qui sont les plus semblables aux prototypes (c'est-à-dire, les véhicules qui sont « les plus proches voisins ») et ainsi déterminer les modèles avec lesquels ils seront en concurrence.

Le constructeur a collecté des données concernant les modèles existants dans plusieurs catégories et a ajouté les détails de ses prototypes. Les catégories dans lesquelles les modèles doivent être comparés comprennent le prix en milliers (*price*), la taille du moteur (*engine_s*), la puissance en chevaux (*horsepow*), l'empattement (*wheelbas*), la largeur (*width*), la longueur (*length*), le poids total (*curb_wgt*), la capacité du réservoir (*fuel_cap*) et le rendement énergétique (*mpg*).

Cet exemple utilise le flux nommé *car_sales_knn.str*, disponible dans le dossier *Demos* du sous-dossier des *flux*. Le fichier de données est *car_sales_knn_mod.sav*. Pour plus d'informations, voir la rubrique «Dossier Demos», à la page 4.

Création du flux



Figure 386. Flux d'échantillons de modélisation KNN

Créez un nouveau flux et ajoutez un noeud source Fichier de statistiques pointant vers *car_sales_knn_mod.sav* dans le dossier *Demos* de votre installation d'IBM SPSS Modeler.

En premier lieu, examinons les données collectées par le constructeur.

- 1. Reliez un noeud Table au noeud source Fichier de statistiques.
- 2. Ouvrez le noeud Table, puis cliquez sur Exécuter.

違 <u>F</u> ile	📄 Edit	🏷 Ger	nerate		9	14	10			0
Table ,	Annotations									
	manufact	model	sales	resale	type	price	engine_s	horsepow	wheelbas	width
140	Toyota	Celica	33.269	15.445	0.0	16	1.800	140.000	102.400	68.3
141	Toyota	Tacoma	84.087	9.575	1.0	11	2.400	142.000	103.300	66.5 1
142	Toyota	Sienna	65.119	\$null\$	1.0	22	3.000	194.000	114.200	73.4 1
143	Toyota	RAV4	25.106	13.325	1.0	16	2.000	127.000	94.900	66.7 1
144	Toyota	4Run	68.411	19.425	1.0	22	2.700	150.000	105.300	66.5 1
145	Toyota	Land	9.835	34.080	1.0	51	4.700	230.000	112.200	76.4 1
146	Volksw	Golf	9.761	11.425	0.0	14	2.000	115.000	98.900	68.3 1
147	Volksw	Jetta	83.721	13.240	0.0	16	2.000	115.000	98.900	68.3 1
148	Volksw	Passat	51.102	16.725	0.0	21	1.800	150.000	106.400	68.5 1
149	Volksw	Cabrio	9.569	16.575	0.0	19	2.000	115.000	97.400	66.7 1
150	Volksw	GTI	5.596	13.760	0.0	17	2.000	115.000	98.900	68.3 1
151	Volksw	Beetle	49.463	\$null\$	0.0	15	2.000	115.000	98.900	67.9 1
152	Volvo	S40	16.957	\$null\$	0.0	23	1.900	160.000	100.500	67.6
153	Volvo	V40	3.545	\$null\$	0.0	24	1.900	160.000	100.500	67.6
154	Volvo	S70	15.245	\$null\$	0.0	27	2.400	168.000	104.900	69.3
155	Volvo	V70	17.531	\$null\$	0.0	28	2.400	168.000	104.900	69.3
156	Volvo	C70	3.493	\$null\$	0.0	45	2.300	236.000	104.900	71.5
157	Volvo	S80	18.969	\$null\$	0.0	36	2.900	201.000	109.900	72.1 :
158		newC	\$null\$	\$null\$	\$n	21	1.500	76.000	106.300	67.9 :
159		newT	\$null\$	\$null\$	\$n	34	3.500	167.000	109.800	75.2
	4						117			

Figure 387. Données source pour les automobiles et les camions

Les détails des deux prototypes, nommés *newCar* et *newTruck*, ont été ajoutés à la fin du fichier. Nous pouvons voir à partir des données source que le constructeur utilise la classification « camion » (valeur 1 dans la colonne *type*) de manière très globale pour signifier tout type de véhicule non automobile. La dernière colonne, *partition*, est nécessaire car les deux prototypes peuvent être désignés comme des ensembles de rétention lorsque nous sommes amenés à identifier leurs voisins les plus proches. De cette manière, leur données n'influencent pas les calculs car il s'agit du reste du marché que nous souhaitons prendre en considération. Le réglage de la valeur de la *partition* des deux enregistrements de rétention sur 1, alors que tous les autres enregistrements ont une valeur de 0 dans ce champ, nous permettra d'utiliser ce champ par la suite, lorsque nous réglerons les enregistrements centraux (les enregistrements pour lesquels nous souhaitons calculer les voisins les plus proches).

Laissons la fenêtre des sorties du tableau ouverte pour le moment, car nous la consulterons par la suite.

Type	review				0	
	Read Vi	alues Cle	ar Values	Clear All V	alues	
Field Measurement		Values	Missing	Check	Role	
👷 norsepow	🞸 conunuous	[55.0,450		NONE	🔳 input	
nteelbas 👘	🔗 Continuous	[92.6,138		None	🕥 Input	
🤣 width	🔗 Continuous	[62.6,79.9]		None	🔪 Input	
🚯 length	🔗 Continuous	[149.4,22		None	🔪 Input	
🔁 curb wat	Continuous	[1.895,5		None	🔪 Input	
fuel cap	Continuous	[10.3.32.0]		None	> Input	
mpa	Continuous	[15.0.46.0]		None		
	Continuous	[-2.20727		None		
partition	Flag	1.0/0.0		None		
 ➢ Insales ➢ partition O View curren 	Continuous Flag t fields O View unu	[-2.20727 1.0/0.0 sed field settin	ngs	None None	None	

Figure 388. Paramètres du noeud type

- 3. Ajoutez un noeud type au flux.
- 4. Reliez un noeud type au noeud source Fichier de statistiques.
- 5. Ouvrez le noeud type.

Nous souhaitons effectuer la comparaison uniquement sur les champs *price* à *mpg*, aussi laissons-nous le rôle de tous ces champs configuré sur **Entrée**.

- 6. Configurez le rôle de tous les autres champs (manufact à type, plus lnsales) sur Aucun.
- 7. Configurez le niveau de mesure du dernier champ, *partition*, sur **Indicateur**. Vérifiez que son rôle est configuré sur **Entrée**.
- 8. Cliquez sur Lire les valeurs pour lire les valeurs de données dans le flux.
- 9. Cliquez sur OK.



Figure 389. Choisir d'identifier les voisins les plus proches

- 10. Reliez un noeud KNN au noeud type.
- 11. Ouvrez le noeud KNN.

Cette fois-ci, nous n'allons pas prédire un champ cible, car nous souhaitons seulement trouver les voisins les plus proches de nos deux prototypes.

- 12. Dans l'onglet Objectifs, sélectionnez Identifier uniquement les voisins les plus proches.
- 13. Cliquez sur l'onglet Paramètres.

😡 No Targets	×
Objectives Fields	Settings Annotations
Settings	
Model	Model name: O Auto O Custom
Neighbors	☑ Use partitioned data
Feature Selection	Build model for each split
Cross-Validation	To select fields manually, choose "Use custom settings" on the Fields tab
Analyze	Partition:
	Splits:
	×
	Vormalize range inputs
	Use case labels
	✓ Identify focal record artition
OK 🕨 Rur	Cancel Apply Reset

Figure 390. Utilisation du champ de partition pour identifier les enregistrements centraux

Nous pouvons maintenant utiliser le champ *partition* pour identifier les enregistrements centraux (les enregistrements pour lesquelles nous souhaitons identifier les voisins les plus proches). En utilisant un champ indicateur, nous nous assurons que les enregistrements dont la valeur de ce champ est configurée sur 1 deviennent nos enregistrements centraux.

Comme nous l'avons vu, les seuls enregistrements qui possèdent une valeur de 1 pour ce champ sont *newCar* et *newTruck*, aussi ces derniers seront-ils nos enregistrements centraux.

- 14. Dans le panneau Modèle de l'onglet Paramètres, cochez la case Identifier un enregistrement central.
- 15. Dans la liste déroulante de ce champ, sélectionnez **partition**.
- 16. Cliquez sur le bouton Exécuter.

Examen des sorties



Figure 391. Fenêtre Visualiseur de modèles

Un nugget de modèle a été créé dans l'espace de travail du flux et dans la palette Modèles. Ouvrez l'un des nuggets pour afficher le Visualiseur de modèles qui dispose d'une fenêtre à deux panneaux :

- Le premier affiche une présentation du modèle, appelée vue principale. La vue principale du modèle Voisin le plus proche est aussi appelée **espace du prédicteur**.
- Le second affiche un des deux types de vues :

Une vue de modèle auxiliaire affiche davantage d'informations sur le modèle, mais n'est pas focalisée sur le modèle lui-même.

Un vue liée est une vue montrant les détails d'une caractéristique du modèle lorsque vous faites défiler une partie de la vue principale.

Espace du prédicteur



This chart is a lower-dimensional projection of the predictor space, which contains a total of 9 predictors.

Figure 392. Graphique de l'espace du prédicteur

Le graphique de l'espace du prédicteur est un graphique 3D interactif qui représente les points des données pour les trois caractéristiques (les trois premiers champs d'entrée des données source) représentant le prix, la taille du moteur et la puissance.

Nos deux enregistrements centraux sont mis en surbrillance en rouge, avec des lignes qui les relient à leurs k voisins les plus proches.

En cliquant sur le graphique et en le faisant glisser, vous pouvez le faire pivoter et obtenir une meilleure vue de la distribution des points dans l'espace du prédicteur. Cliquez sur le bouton **Réinitialiser** pour rétablir la vue par défaut.

Graphique des homologues



Figure 393. Graphique des homologues

La vue auxiliaire par défaut est le graphique des homologues qui met en évidence les deux enregistrements centraux sélectionnés dans l'espace du prédicteur ainsi que leurs *k* voisins les plus proches pour chacune des six caractéristiques (les six premiers champs d'entrée des données source).

Les véhicules sont représentés par leur numéro d'enregistrement dans les données source. Nous devons maintenant établir les sorties à partir du noeud Table afin de les identifier.

Si la sortie du noeud Table est encore disponible :

- 1. Cliquez sur l'onglet **Sorties** du panneau du gestionnaire, en haut à droite de la fenêtre principale d'IBM SPSS Modeler.
- 2. Double-cliquez sur l'entrée **Table (16 champs, 159 enregistrements)**. Si la sortie de la table n'est plus disponible :
- 3. Dans la fenêtre principale d'IBM SPSS Modeler, cliquez sur le noeud Table.
- 4. Cliquez sur Exécuter.

違 <u>F</u> ile	📄 Edit	🕙 Ger	nerate			14	11			0
Table ,	Annotations									
	manufact	model	sales	resale	type	price	engine_s	horsepow	wheelbas	width
140	Toyota	Celica	33.269	15.445	0.0	16	1.800	140.000	102.400	68.3
141	Toyota	Tacoma	84.087	9.575	1.0	11	2.400	142.000	103.300	66.5
142	Toyota	Sienna	65.119	\$null\$	1.0	22	3.000	194.000	114.200	73.4
143	Toyota	RAV4	25.106	13.325	1.0	16	2.000	127.000	94.900	66.7 :
144	Toyota	4Run	68.411	19.425	1.0	22	2.700	150.000	105.300	66.5 1
145	Toyota	Land	9.835	34.080	1.0	51	4.700	230.000	112.200	76.4
146	Volksw	Golf	9.761	11.425	0.0	14	2.000	115.000	98.900	68.3 1
147	Volksw	Jetta	83.721	13.240	0.0	16	2.000	115.000	98.900	68.3 1
148	Volksw	Passat	51.102	16.725	0.0	21	1.800	150.000	106.400	68.5 1
149	Volksw	Cabrio	9.569	16.575	0.0	19	2.000	115.000	97.400	66.7 :
150	Volksw	GTI	5.596	13.760	0.0	17	2.000	115.000	98.900	68.3 1
151	Volksw	Beetle	49.463	\$null\$	0.0	15	2.000	115.000	98.900	67.9 1
152	Volvo	S40	16.957	\$null\$	0.0	23	1.900	160.000	100.500	67.6 1
153	Volvo	V40	3.545	\$null\$	0.0	24	1.900	160.000	100.500	67.6 1
154	Volvo	S70	15.245	\$null\$	0.0	27	2.400	168.000	104.900	69.3
155	Volvo	V70	17.531	\$null\$	0.0	28	2.400	168.000	104.900	69.3
156	Volvo	C70	3.493	\$null\$	0.0	45	2.300	236.000	104.900	71.5
157	Volvo	S80	18.969	\$null\$	0.0	36	2.900	201.000	109.900	72.1
158		newC	\$null\$	\$null\$	\$n	21	1.500	76.000	106.300	67.9
159	100	newT	\$null\$	\$null\$	\$n	34	3.500	167.000	109.800	75.2
	4						1111			

Figure 394. Identification des enregistrements par numéro d'enregistrement

Si vous accédez en bas de la table, vous pouvez constater que *newCar* et *newTruck* sont les deux derniers enregistrements des données avec les numéros 158 et 159, respectivement.



Figure 395. Comparaison des caractéristiques dans le graphique des homologues

A partir de ceci, nous pouvons voir dans le graphique des homologues, par exemple, que *newTruck* (159) possède une plus grande taille de moteur que tous ses voisins les plus proches, alors que *newCar* (158) possède un moteur plus petit que tous *ses* voisins les plus proches.

Pour chacune des six caractéristiques, vous pouvez déplacer la souris sur les points individuels afin d'afficher la valeur réelle de chacune des caractéristiques de cette observation spécifique.

Mais quels sont les véhicules qui sont les voisins les plus proches de *newCar* et de *newTruck* ? Comme le graphique des homologues est un peu encombré, simplifions la vue.

- 5. Cliquez sur la liste déroulante Vue en bas du graphique des homologues (l'entrée nommée Pairs).
- 6. Sélectionnez Tableau des voisins et des distances.

Tableau des voisins et des distances

	Displa	yed for Ir	nitial Foc	al Records				
- ID I	Nea	rest Neigh	bors	Nea	rest Distar			
Focal Record	1	2	3	1	2			
158	131	130	58	0.979	0.990			
159	105	92	101	0.580	0.634			

k Nearest Neighbors and Distances

Figure 396. Tableau des voisins et des distances

Ceci est mieux. Nous pouvons maintenant voir les trois modèles dont nos deux prototypes sont les voisins les plus proches sur le marché.

Pour *newCar* (enregistrement central 158) il s'agit de la Saturn SC (131), de la Saturn SL (130) et de la Honda Civic (58).

Ceci n'est pas vraiment surprenant (les trois sont des berlines de taille moyenne, donc *newCar* correspond bien, en particulier en ce qui concerne son excellente efficacité énergétique).

Pour *newTruck* (enregistrement central 159), les voisins les plus proches seront la Nissan Quest (105), la Mercury Villager (92) et la Mercedes M-Class (101).

Comme nous l'avons vu plus tôt, il ne s'agit pas nécessairement de camions au sens propre, mais simplement de véhicules qui ne sont pas classés comme des automobiles. Si nous observons la sortie du noeud Table afin de rechercher les voisins les plus proches, nous constatons que *newTruck* est relativement cher tout en étant l'un des plus lourds de sa catégorie. Cependant, l'efficacité énergétique est une fois encore meilleure que ses plus proches rivaux, ce qui devrait jouer en sa faveur.

Récapitulatif

Nous avons vu comment vous pouvez utiliser l'analyse des voisins les plus proches pour comparer un ensemble étendu de caractéristiques d'observations à partir d'un jeu de données particulier. Nous avons aussi calculé, pour deux enregistrements de rétention très différents, les observations qui ressemblent le plus étroitement à ces ensembles de rétention.

Chapitre 29. Découverte des relations de causalité dans les métriques métier (TCM)

Une entreprise suit divers indicateurs clé de performance qui décrivent l'état financier de l'entreprise au fil du temps, ainsi qu'un nombre de métriques qu'elle peut contrôler. Elle veut utiliser la modélisation de causalité temporelle pour découvrir les relations de causalité entre les métriques contrôlables et les indicateurs clé de performance. Elle veut également connaître les relations de causalité entre les indicateurs clé de performance.

Le fichier de données tcm_kpi.sav contient les données hebdomadaires des indicateurs clé de performance et les métriques contrôlables. Les données des indicateurs clé de performance sont stockées dans des zones ayant le préfixe *KPI*. Les données des métriques contrôlables sont stockées dans des zones ayant le préfixe *Lever*.

Création du flux



Figure 397. Exemple de flux de modélisation TCM

- 1. Créez un flux et ajoutez un noeud source de fichier de statistiques au fichier *tcm_kpi.sav* dans le dossier *Demos* de l'installation IBM SPSS Modeler.
- 2. Reliez un noeud Table au noeud source Fichier de statistiques.
- **3.** Ouvrez le noeud Table et cliquez sur **Exécuter** pour consulter les données. Il contient les données hebdomadaires des indicateurs clé de performance et les métriques contrôlables. Les données des indicateurs clé de performance sont stockées dans des zones ayant le préfixe *KPI*, et les données des métriques contrôlables sont stockées dans des zones ayant le préfixe *Lever*.

🔟 Table	(31 fields, 112 re	cords)							x
じ <u>F</u> ile	📄 <u>E</u> dit 🛛 💐) <u>G</u> enerat	e 健		***			0	×
Table	Annotations								
	date	Lever1	Lever2	Lever3	Lever4	Lever5	KPI_1	KPI_2	
1	2008-09-07	6.816	1.176	101.839	88.258	2027.711	1.829	1891.833	-
2	2008-09-14	6.091	1.172	120.610	103.803	2343.404	2.162	2125.261	
3	2008-09-21	8.108	1.093	70.512	81.053	1813.224	1.809	1848.765	
4	2008-09-28	6.503	1.121	78.581	86.393	2722.012	1.784	2551.153	
5	2008-10-05	8.564	1.024	148.985	104.379	2235.634	1.704	2186.098	
6	2008-10-12	7.331	0.848	170.236	91.477	2607.424	1.642	1711.295	
7	2008-10-19	6.996	1.362	239.189	69.636	2354.322	1.681	2112.309	
8	2008-10-26	7.863	0.959	169.925	87.400	1860.496	2.304	1561.226	
9	2008-11-02	7.894	1.131	307.334	109.800	1600.156	1.782	1929.897	
10	2008-11-09	6.548	1.052	467.642	77.574	2007.203	1.913	2042.415	
11	2008-11-16	4.281	1.232	564.812	80.350	1764.707	1.915	2268.544	
12	2008-11-23	7.458	1.219	523.018	105.373	2106.771	1.676	2451.158	
13	2008-11-30	7.235	0.978	628.724	73.206	2666.294	2.160	2558.336	
14	2008-12-07	7.752	1.032	654.648	99.905	1915.698	1.964	1614.402	
15	2008-12-14	7.839	0.770	712.274	80.301	1811.261	1.147	1925.271	
16	2008-12-21	8.529	1.374	699.621	98.391	1792.807	2.033	2320.790	
17	2008-12-28	6.069	1.034	562.279	117.396	2216.657	0.879	2478.630	
18	2009-01-04	6.174	1.442	613.071	72.062	2530.900	1.701	1769.694	
19	2009-01-11	7.046	1.410	718.218	95.594	2285.149	1.841	2215.692	
20	2009-01-18	5.805	0.933	908.362	83.863	2391.528	1.977	2094.555	-
	4							1	
								C	ж

Figure 398. Données source des indicateurs clés de performance et des métriques contrôlables

- 4. Ajoutez un noeud type au flux.
- 5. Reliez un noeud type au noeud source Fichier de statistiques.

Exécution d'une analyse

1. association d'un noeud TCM au noeud Type, puis ouvrez le noeud TCM et accédez à la section **Observations** de l'onglet **Zones**.

🕐 ТСМ	
9.	
Fields Data Spec	ifications Build Options Model Options Annotations
Select an item.	
Observations	Observations are specified by a date/time field
Observations	Date/time field:
	🔏 date 🚽
	Time interval: Weeks
	Observations are defined as periods or cyclic periods
	P <u>e</u> riod field:
	(none) Increment by: 1 🖨 Starting value: 1 🖨
	Cycle fields:
	Level Field Cycle length Starting value
	↓
OK 🕨 Rur	Cancel Apply Reset

Figure 399. Modélisation de causalité temporelle, boîte de dialogue d'introduction

- 2. Sélectionnez date dans la zone de date/heure et Semaines dans la zone d'intervalle de temps.
- 3. Cliquez sur Séries temporelles et sélectionnez Utiliser des rôles prédéfinis.

Dans le jeu de données d'échantillon tcm_kpi.sav, les zones *Lever1* à *Lever5* ont le rôle d'entrée, et *KPI_1* à *KPI_25* ont le rôle Les deux. Lorsque vous sélectionnez **Utiliser des rôles prédéfinis**, les zones ayant un rôle d'entrée sont traitées comme entrées possibles, et les zones ayant le rôle Les deux sont traitées comme entrées de la modélisation de causalité temporelle.

La procédure de modélisation de causalité temporelle détermine les meilleures entrées pour chaque cible depuis l'ensemble d'entrées possibles. Dans cet exemple, les entrées possibles sont les zones *Lever1* à *Lever5* et les zones *KPI_1* à *KPI_25*.

4. Cliquez sur Exécuter.

Graphique de qualité du modèle de système global

L'élément de sortie Qualité du modèle de système global, généré par défaut, affiche un diagramme à barres et le tracé de points associés de l'ajustement de modèle pour tous les modèles. Il existe un modèle distinct pour chaque série cible. L'ajustement de modèle est mesuré par la statistique d'ajustement choisie. Cet exemple utilise la statistique d'ajustement par défaut, à savoir R-carré.

L'élément Qualité du modèle global contient des fonctions interactives. Pour les activer, activez l'élément en cliquant deux fois sur le graphique Système de modèle global dans le visualiseur.



Figure 400. Qualité globale du modèle

Cliquez sur une barre dans le diagramme à barres pour filtrer le tracé de points et afficher uniquement les modèles associés à la barre sélectionnée. Placez le pointeur de la souris sur un point pour afficher une infobulle qui contient le nom des séries associées et la valeur de la statistique d'ajustement. Vous pouvez rechercher le modèle d'une série cible dans le tracé de points en définissant le nom de la série dans la zone **Rechercher un modèle pour la cible**.

Système de modèle global

L'élément de sortie Système de modèle global, généré par défaut, affiche la représentation graphique des relations de causalité entre les séries dans le système de modèle. Par défaut, les relations des 10 principaux modèles sont indiquées, telles que déterminées par la valeur de la statistique d'ajustement R-carré. Le nombre de modèles principaux (appelés également modèles de meilleur ajustement) et la statistique d'ajustement sont définis dans les paramètres de série à afficher (dans l'onglet Options de création) de la boîte de dialogue Modélisation de causalité temporelle.

L'élément Système de modèle global contient des fonctions interactives. Pour les activer, activez l'élément en cliquant deux fois sur le graphique Système de modèle global dans le visualiseur. Dans cet exemple, il est très important de voir les relations entre toutes les séries dans le système. Dans la sortie interactive, sélectionnez **Toutes les séries** dans la liste déroulante **Mettre en évidence les relations de**.



Figure 401. Système de modèle global, vue de toutes les séries

Toutes les lignes qui connectent une cible à ses sorties ont la même couleur, et la flèche sur chaque ligne pointe d'une entrée vers la cible de l'entrée. Par exemple, *Lever3* est une entrée vers *KPI_19*.

L'épaisseur de chaque ligne indique la signification de la relation de causalité, les lignes épaisses indiquant une relation plus significative. Par défaut, les relations de causalité avec une valeur de signification supérieure à 0,05 sont masquées. Au niveau 0,05, seuls *Lever1*, *Lever3*, *Lever4* et *Lever5* ont des relations de causalité significatives avec les zones d'indicateur clé de performance. Vous pouvez changer le niveau d'importance seuil en entrant une valeur dans la zone **Masquer les liens avec une valeur de signification supérieure à**.

Outre les relations de causalité entre les zones *Lever* et les zones d'indicateur clé de performance, l'analyse a découvert les relations entre les zones d'indicateur clé de performance. Par exemple, *KPI_10* a été sélectionné comme entrée dans le modèle pour *KPI_2*.

Vous pouvez filtrer la vue pour afficher uniquement les relations d'une seule série. Par exemple, pour afficher uniquement les relations *KPI_19*, cliquez sur le libellé de *KPI_19*, cliquez avec le bouton droit et sélectionnez **Mettre en évidence les relations des séries**.



Figure 402. Système de modèle global, vue d'une seule série

Cette vue contient les entrées de *KPI_19* ayant une valeur significative inférieure ou égale à 0,05. Elle indique également qu'au niveau d'importance 0,05, *KPI_19* a été sélectionné comme entrée de *KPI_18* et *KPI_7*.

Outre l'affichage des relations de la série sélectionnée, l'élément de sortie contient également des informations sur les valeurs extrêmes détectées pour la série. Cliquez sur l'onglet **Séries avec valeurs extrêmes**.

Series N	with	Outliers	for	KPI_	19
----------	------	----------	-----	------	----

Series	Time	Observed Value
KPI_19	2008-10-12	7,358,201.68
	2009-04-05	2.10E+007
	2010-09-19	6,492,157.97

Figure 403. Valeurs extrêmes pour KPI_19

Trois valeurs extrêmes ont été détectées pour *KPI_19*. Compte tenu du système de modèle, qui contient toutes les connexions découvertes, il est possible d'aller au-delà de la détection des valeurs extrêmes, et
de déterminer la série qui génère probablement une valeur extrême. Ce type d'analyse s'appelle une analyse de la cause première de valeur extrême ; elle est traitée dans une rubrique ultérieure dans cette étude de cas.

Diagrammes d'impact

Vous pouvez obtenir une vue complète de toutes les relations associées à une série en générant un diagramme d'impact. Cliquez sur le libellé de *KPI_19* dans le graphique Système de modèle global, cliquez avec le bouton droit et sélectionnez **Créer un diagramme d'impact**.



Figure 404. Diagramme d'impact des effets

Lorsqu'un diagramme d'impact est créé depuis le système de modèle global, comme dans cet exemple, il contient initialement les séries affectées par la série sélectionnée. Par défaut, les diagrammes d'impact contiennent trois niveaux d'effet, le premier niveau étant la série d'intérêt. Chaque niveau supplémentaire indique des effets indirects supplémentaires de la série d'intérêt. Vous pouvez changer le **nombre de niveaux à afficher** pour afficher plus ou moins de niveaux d'effet. Le diagramme d'impact de cet exemple montre que *KPI_19* est une entrée directe de *KPI_18* et *KPI_7*, et qu'il affecte indirectement des séries du fait de son effet sur la série *KPI_7*. Comme dans le système de modèle global, l'épaisseur des lignes indique la signification des relations de causalité.

Le graphique affiché dans chaque noeud du diagramme d'impact indique les dernières valeurs L+1 des séries associées à la fin de la période d'estimation, et les valeurs prévisionnelles, où L est le nombre de termes de décalage inclus dans chaque modèle. Vous pouvez obtenir un graphique séquentiel de ces valeurs en cliquant sur le noeud associé.

Lorsque vous cliquez deux fois sur un noeud, vous définissez les séries associées comme séries d'intérêt et régénérez le diagramme d'impact en fonction de la série. Vous pouvez également définir un nom de série dans la zone de **série d'intérêt** pour sélectionner une série d'intérêt différente.

Les diagrammes d'impact peuvent également afficher les séries qui affectent la série d'intérêt. Ces séries s'appellent des *causes*. Pour afficher les séries qui affectent *KPI_19*, sélectionnez **Causes des séries** dans la liste déroulante **Afficher**.



Figure 405. Diagramme d'impact des causes

Cette vue indique que le modèle pour *KPI_19* a quatre entrées et que *Lever3* a la connexion de causalité la plus significative avec *KPI_19*. Il indique également les séries qui affectent indirectement *KPI_19* du fait de leurs effets sur *KPI_7* et *KPI_17*. Le même concept de niveaux abordé pour les effets s'applique aux causes. De même, vous pouvez changer le **nombre de niveaux à afficher** pour afficher plus ou moins de causes.

Détermination des causes premières des valeurs extrêmes

Avec un système de modèle, il est possible d'aller au-delà de la détection des valeurs extrêmes, et de déterminer la série qui génère probablement une valeur extrême. Ce processus s'appelle une analyse de la cause première de valeur extrême et il doit être exécuté en fonction de chaque série. L'analyse nécessite un système de modèle de causalité temporelle et les données qui ont été utilisées pour créer le système. Dans cet exemple, le jeu de données actif est les données utilisées pour créer le système de modèle.

Pour exécuter une analyse de cause première :

1. Dans la boîte de dialogue TCM, accédez à l'onglet **Options de création** et cliquez sur **Série à afficher** dans la liste **Sélectionner un élément**.

🕐 ТСМ				X		
90				0		
Fields Data Spe	ifications Build Options	Model Options	Annotations			
<u>S</u> elect an item:	_					
General	Display targets associ	ated with best-fitting	g models			
Series to Display	Fixed number of	targets				
Output Options	<u>N</u> umber: 10					
Estimation Period	Percentage of total number of targets					
	Percenta <u>a</u> e:					
	Goo <u>d</u> ness of fit mea	asure: R square		~		
	Specify Individual Serie	es				
	<u>F</u> ields:			Fields to Displa <u>y</u> :		
	Sort: None	-		KPI_19		
	KPI_11	4		•		
	✓ KPI_12 ✓ KPI_13		•	The second secon		
	KPI_14					
	KPI_15	-				
	All 🖋					
OK 🕨 R	Cancel			<u>Apply</u> <u>R</u> eset		

Figure 406. Séries à afficher du modèle de causalité

- 2. Transférez *KPI_19* vers la liste **Zones à afficher**.
- 3. Cliquez sur **Options de sortie** dans la liste **Sélectionner un élément** dans l'onglet Options.

📀 тсм		— ×-
90		0 - 1
Fields Data Specif	fications Build Options Model Options Annotations	
<u>S</u> elect an item:		
General	Output for targets	
Series to Display	Overall model system	Series plot
Output Options	Significance level: 0.05	Residuals plot
Estimation Period	Model fit statistics and outliers	Top Inputs
	Model effects and model parameters	Eorecast table
	Impact diagram	
	Type: Both causes and effects	
	Output for series	
	Same as for targets	
	Overall model system	Series plot
	Significance level: 0.05	Residuals plot
	Model fit statistics and outliers	Top Inputs
	Model effects and model parameters	Eorecast table
	Impact diagram	
	Type: Both causes and effects	
	✓ Outlier root cause analysis	
	Model fit across all models	me
	R sguare BIC Series transfor	rmations
	Root mean square per <u>c</u> entage error AIC	
	Root mean square error	
OK Run	Cancel	Apply Reset

Figure 407. Options de sortie du modèle de causalité temporelle

- 4. Désélectionnez Système de modèle global, Comme pour les cibles, R carré et Transformations des séries.
- 5. Sélectionnez **Analyse de la cause première de valeur extrême** et conservez les paramètres existants pour **Sortie** et **Niveaux de causalité**.
- 6. Cliquez sur Exécuter.
- 7. Cliquez deux fois sur le graphique Analyse de la cause première de valeur extrême pour *KPI_19* dans le visualiseur pour l'activer.



Figure 408. Analyse de la cause première de valeur extrême pour KPI_19

Les résultats de l'analyse sont récapitulés dans le tableau Valeurs extrêmes. Le tableau indique que des causes premières ont été trouvées pour les valeurs extrêmes 2009-04-05 et 2010-09-19, mais pas pour la valeur extrême 2008-10-12. Cliquez sur une ligne dans le tableau des valeurs extrêmes pour mettre en évidence la série de la cause première, comme indiqué ici pour la valeur extrême 2009-04-05. Cette action met également en évidence la valeur extrême sélectionnée dans le graphique séquentiel. Vous pouvez également cliquer sur une valeur extrême directement pour mettre en évidence le chemin de la série de cause première de la valeur.

Pour la valeur extrême 2009-04-05, la cause première est *Lever3*. Le diagramme montre que *Lever3* est l'entrée directe de *KPI_19* et qu'il impacte indirectement *KPI_19* par ses effets sur les autres séries qui affectent *KPI_19*. L'un des paramètres configurables de l'analyse de la cause première de valeur extrême est le nombre de niveaux de causalité utilisés pour rechercher les causes premières. Par défaut, la recherche s'effectue dans trois niveaux. Les occurrences de la série de cause première s'affichent jusqu'au nombre défini de niveaux de causalité. Dans cet exemple, *Lever3* existe au premier et au troisième niveaux de cause première.

Chaque noeud dans le chemin en évidence d'une valeur extrême contient un graphique dont la plage de temps dépend du niveau du noeud. Pour les noeuds du premier niveau de cause première, la plage est T-1 à T-L, où T est le moment d'occurrence de la valeur extrême et L est le nombre de termes de décalage inclus dans chaque modèle. Pour les noeuds dans le deuxième niveau de causalité, la plage est T-2 à T-L-1, et pour le troisième niveau, la plage est T-3 à T-L-2. Vous pouvez obtenir un graphique séquentiel détaillé de ces valeurs en cliquant sur le noeud associé.

Exécution de scénario

Avec un système de modèle de causalité temporelle, vous pouvez exécuter des scénarios définis par l'utilisateur. Un *scénario* est défini par une série temporelle appelée *série racine* et définit un ensemble de valeurs spécifiées par l'utilisateur pour la série sur une période donnée. Les valeurs définies sont utilisées pour générer les prévisions pour les séries temporelles affectées par les séries racine. L'analyse nécessite un système de modèle de causalité temporelle et les données qui ont été utilisées pour créer le système. Dans cet exemple, le jeu de données actif est les données utilisées pour créer le système de modèle.

Pour exécuter les scénarios :

- 1. Dans la boîte de dialogue de sortie TCM, cliquez sur le bouton Analyse de scénario.
- 2. Dans la boîte de dialogue des scénarios du modèle de causalité temporelle, cliquez sur **Définir la période du scénario**.

	Date				
Start	2008-09-07				
End	nd 2010-10-24				
Time interval: Weeks					
e Derind for Scenarios					
Specify by start, and and pred	ict through times				
Start of scenario values	Date vvv-MM-dd				
End of scenario values	vvv-MM-dd				
Predict through	vvv-MM-dd				

Figure 409. Période du scénario

- 3. Sélectionnez Spécifier par intervalles de temps relatifs à la fin de la période d'estimation.
- 4. Entrez -3 pour l'intervalle de début et 0 pour l'intervalle de fin.

Ces paramètres indiquent que chaque scénario repose sur des valeurs définies pour les quatre derniers intervalles dans la période d'estimation. Pour cet exemple, les quatre derniers intervalles signifie les quatre dernières semaines. La période sur laquelle les valeurs du scénario sont définies s'appelle la *période du scénario*.

- Entrez 4 pour les intervalles pour effectuer des prévisions après la fin des valeurs du scénario.
 Ce paramètre indique que les prévisions sont générées pour quatre intervalles après la fin de la période du scénario.
- 6. Cliquez sur **Continuer**
- 7. Cliquez sur Ajouter un scénario dans l'onglet Scénarios.

			Rgot field: Provide the currently defined targets By default, affected targets up to the currently defined maximum of 25 are automatically determined. Affected targets::
Scenario Definition - Scenario ID:	ever3_25pct alues are applied to th	ne data used for mod	deling, after any aggregation or distribution of the original data.
Scenario Definition Scenario ID:	ever3_25pct alues are applied to th ario <u>v</u> alues for root field Date	e data used for mod	deling, after any aggregation or distribution of the original data.
Scenario Definition - Scenario ID: L Scenario v Scenario v © Specify Scena Interval -3	ever3_25pct alues are applied to th ario <u>v</u> alues for root field Date 2010-10-03	e data used for mod d Scenario value	deling, after any aggregation or distribution of the original data.
Scenario Definition - Scenario ID: L Scenario v Scenario v © Specify Scena Interval -3 -2	ever3_25pct alues are applied to th ario values for root field Date 2010-10-03 2010-10-10	e data used for mod	deling, after any aggregation or distribution of the original data. Root field value <read> <read></read></read>
Scenario Definition - Scenario ID: L Scenario v Scenario v © Specify Scena -3 -2 -1 -2	ever3_25pct alues are applied to th ario values for root field Date 2010-10-03 2010-10-10 2010-10-17	e data used for mod	deling, after any aggregation or distribution of the original data. Root field value <read> <read> <read></read></read></read>
Scenario Definition - Scenario ID: L Scenario v Scenario v © Specify Scenario -3 -2 -1 0	ever3_25pct alues are applied to th ario values for root field 2010-10-03 2010-10-10 2010-10-17 2010-10-24	e data used for mod	deling, after any aggregation or distribution of the original data. Root field value <read> <read> <read> <read></read></read></read></read>

Figure 410. Définition de scénario

- 8. Transférez *Lever3* vers la zone **Zone racine** pour déterminer comment les valeurs définies de *Lever3* dans la période du scénario affectent les prévisions des autres séries affectées par causalité par *Lever3*.
- 9. Entrez Lever3_25pct pour l'ID de scénario.
- 10. Sélectionnez **Spécifier une expression pour les valeurs de scénario pour la zone racine** et entrez Lever3*1.25 pour l'expression.

Ce paramètre indique que les valeurs de *Lever3* dans la période de scénario sont 25 % plus élevées que les valeurs observées. Pour plus d'informations sur les expressions complexes, vous pouvez utiliser le générateur d'expression en cliquant sur l'icône de calculatrice.

- 11. Cliquez sur **Continuer**
- 12. Répétez les étapes 10 à 14 pour définir un scénario ayant *Lever3* comme zone racine, Lever3_50pct comme ID de scénario et Lever3*1.5 comme expression.

emporal Causal Model Scenarios			— ×
Options			
Scenario Period			
the time period over which	scenarios are carried out must be defined before scenari	os can be created.	
Define Scenario Period			
Scenarios:			
Scenario ID	Root field	Scenario values	
Lever3 25pct	Lever3	Lever3*1.25	X
Lever3_50pct	Lever3	Lever3*1.5	
Add Scenario Edit Scenari	0		
Add Scenario Edit Scenari	0		
Add Scenario Edit Scenari			

Figure 411. Scénarios

- 13. Cliquez sur l'onglet Options et entrez 2 comme niveau maximum pour les cibles affectées.
- 14. Cliquez sur Exécuter.
- 15. Cliquez deux fois sur le diagramme d'impact de Lever3_50pct dans le visualiseur pour le visualiser.



Figure 412. Diagramme d'impact du scénario : Lever3_50pct

Le diagramme d'impact montre les séries affectées par la série racine *Lever3*. Deux niveaux d'effets sont indiqués, car vous avez défini 2 comme niveau maximum pour les cibles affectées.

Le tableau des valeurs prévues contient les prévisions pour toutes les séries affectées par *Lever3* jusqu'au second niveau des effets. Les prévisions des séries cible dans le premier niveau d'effets commencent à la première période après le début de la période de scénario. Dans cet exemple, les prévisions pour les séries cible dans le premier niveau commencent à 2010-10-10. Les prévisions pour les séries cible dans le second niveau des effets commencent à la seconde période après le début de la période de scénario. Dans cet exemple, les prévisions pour les séries cible dans le second niveau des effets commencent à la seconde période après le début de la période de scénario. Dans cet exemple, les prévisions pour les séries cible dans le second niveau commencent à 2010-10-17. Le caractère étalé des prévisions reflète le fait que les modèles de séries temporelles reposent sur des valeurs de décalage.

16. Cliquez sur le noeud pour que KPI_5 génère un diagramme séquentiel détaillé.



Figure 413. Graphique séquentiel pour KPI_5

Le graphique séquentiel montrent les valeurs prédéfinies du scénario, ainsi que les valeurs des séries en l'absence de scénario. Lorsque la période de scénario contient des périodes dans la période d'estimation, les valeurs observées des séries s'affichent. Pour les périodes après la fin de la période d'estimation, les prévisions d'origine sont affichées.

Remarques

Ces informations ont été développées pour les produits et services offerts en France. Elles peuvent être disponibles dans d'autres langues auprès d'IBM. Toutefois, une copie du produit ou de la version du produit dans cette langue peut être nécessaire pour y accéder.

Le présent document peut contenir des informations ou des références concernant certains produits, logiciels ou services IBM non annoncés dans ce pays. Pour plus de détails, référez-vous aux documents d'annonce disponibles dans votre pays, ou adressez-vous à votre partenaire commercial IBM. Toute référence à un produit, programme ou service IBM n'implique pas que seul ce produit, programme ou service IBM puisse être utilisé. Tout produit, programme ou service fonctionnellement équivalent peut être utilisé s'il n'enfreint aucun droit de propriété intellectuelle d'IBM. Cependant l'utilisateur doit évaluer et vérifier l'utilisation d'un produit, programme ou service non IBM.

IBM peut détenir des brevets ou des demandes de brevet couvrant les produits mentionnés dans le présent document. L'octroi de ce document n'équivaut aucunement à celui d'une licence pour ces brevets. Vous pouvez envoyer par écrit des questions concernant la licence à :

IBM Director of Licensing IBM Corporation North Castle Drive, MD-NC119 Armonk, NY 10504-1785 US

Pour toute demande au sujet des licences concernant les jeux de caractères codés sur deux octets (DBCS), contactez le service Propriété intellectuelle IBM de votre pays ou adressez vos questions par écrit à :

Intellectual Property Licensing Legal and Intellectual Property Law IBM Japan Ltd. 19-21, Nihonbashi-Hakozakicho, Chuo-ku Tokyo 103-8510, Japan

LE PRESENT DOCUMENT EST LIVRE "EN L'ETAT". IBM DECLINE TOUTE RESPONSABILITE, EXPLICITE OU IMPLICITE, RELATIVE AUX INFORMATIONS QUI Y SONT CONTENUES, Y COMPRIS EN CE QUI CONCERNE LES GARANTIES DE VALEUR MARCHANDE OU D'ADAPTATION A VOS BESOINS. Certaines juridictions n'autorisent pas l'exclusion de garanties explicites ou implicites lors de certaines transactions, par conséquent, il est possible que cet énoncé ne vous concerne pas.

Ces informations peuvent contenir des erreurs techniques ou des erreurs typographiques. Ce document est mis à jour périodiquement. Chaque nouvelle édition inclut les mises à jour. IBM peut, à tout moment et sans préavis, modifier les produits et logiciels décrits dans ce document.

Les références à des sites Web non IBM contenues dans le présent document sont fournies à titre d'information uniquement et n'impliquent en aucun cas une adhésion aux données qu'ils contiennent. Les éléments figurant sur ces sites Web ne font pas partie des éléments du présent produit IBM et l'utilisation de ces sites relève de votre seule responsabilité.

IBM pourra utiliser ou diffuser, de toute manière qu'elle jugera appropriée et sans aucune obligation à votre égard, tout ou partie des informations qui lui seront fournies.

Les licenciés souhaitant obtenir des informations permettant : (i) l'échange des données entre des logiciels créés de façon indépendante et d'autres logiciels (dont celui-ci), et (ii) l'utilisation mutuelle des données ainsi échangées, doivent adresser leur demande à :

IBM Director of Licensing IBM Corporation North Castle Drive, MD-NC119 Armonk, NY 10504-1785 US

Ces informations peuvent être disponibles, soumises à des conditions générales, et dans certains cas payantes.

Le programme sous licence décrit dans le présent document et tous les éléments sous licence disponibles s'y rapportant sont fournis par IBM conformément aux dispositions du Livret Contractuel IBM, des Conditions internationales d'utilisation des Logiciels IBM ou de tout autre contrat équivalent.

Les données de performances et les exemples de client ne sont présentés qu'à des fins d'illustration. Les performances réelles peuvent varier en fonction des configurations et des conditions d'exploitation spécifiques.

Les informations concernant des produits non IBM ont été obtenues auprès des fournisseurs de ces produits, par l'intermédiaire d'annonces publiques ou via d'autres sources disponibles. IBM n'a pas testé ces produits et ne peut confirmer l'exactitude de leurs performances ni leur compatibilité. Aucune réclamation relative à des produits non IBM ne pourra être reçue par IBM. Toute question concernant les performances de produits non IBM doit être adressée aux fournisseurs de ces produits.

Toute instruction relative aux intentions d'IBM pour ses opérations à venir est susceptible d'être modifiée ou annulée sans préavis, et doit être considérée uniquement comme un objectif.

Ces informations contiennent des exemples de données et de rapports utilisés au cours d'opérations quotidiennes standard. Pour les illustrer le mieux possible, ces exemples contiennent des noms d'individus, d'entreprises, de marques et de produits. Toute ressemblance avec des noms de personnes ou de sociétés réelles serait purement fortuite.

Marques

IBM, le logo IBM et ibm.com sont des marques d'International Business Machines dans de nombreux pays. Les autres noms de produits et de services peuvent être des marques d'IBM ou d'autres sociétés. La liste actualisée de toutes les marques d'IBM est disponible sur la page Web "Copyright and trademark information", à l'adresse www.ibm.com/legal/copytrade.shtml.

Adobe, le logo Adobe, PostScript et le logo PostScript sont des marques déposées ou des marques commerciales de Adobe Systems Incorporated aux Etats-Unis et/ou dans d'autres pays.

Intel, le logo Intel, Intel Inside, le logo Intel Inside, Intel Centrino, le logo Intel Centrino, Celeron, Intel Xeon, Intel SpeedStep, Itanium, et Pentium sont des marques commerciales ou des marques déposées de Intel Corporation ou de ses filiales aux Etats-Unis et dans d'autres pays.

Linux est une marque déposée de Linus Torvalds aux Etats-Unis et/ou dans d'autres pays.

Microsoft, Windows, Windows NT et le logo Windows sont des marques commerciales de Microsoft Corporation aux Etats-Unis et/ou dans d'autres pays.

UNIX est une marque déposée de The Open Group aux Etats-Unis et dans d'autres pays.

Les marques commerciales Java et basées sur Java ainsi que les logos sont des marques commerciales ou déposées de Oracle et/ou de ses filiales.

Dispositions relatives à la documentation du produit

Les droits d'utilisation relatifs à ces publications sont soumis aux dispositions suivantes.

Applicabilité

Les présentes dispositions viennent s'ajouter à toute autre condition d'utilisation applicable au site Web IBM.

Usage personnel

Vous pouvez reproduire ces publications pour votre usage personnel, non commercial, sous réserve que toutes les mentions de propriété soient conservées. Vous ne pouvez distribuer ou publier tout ou partie de ces publications ou en faire des oeuvres dérivées sans le consentement exprès d'IBM.

Usage commercial

Vous pouvez reproduire, distribuer et publier ces publications uniquement au sein de votre entreprise, sous réserve que toutes les mentions de propriété soient conservées. Vous ne pouvez reproduire, distribuer, afficher ou publier tout ou partie de ces publications en dehors de votre entreprise, ou en faire des oeuvres dérivées, sans le consentement exprès d'IBM.

Droits

Excepté les droits d'utilisation expressément accordés dans ce document, aucun autre droit, licence ou autorisation, implicite ou explicite, n'est accordé pour ces publications ou autres informations, données, logiciels ou droits de propriété intellectuelle contenus dans ces publications.

IBM se réserve le droit de retirer les autorisations accordées ici si, à sa discrétion, l'utilisation des publications s'avère préjudiciable à ses intérêts ou que, selon son appréciation, les instructions n'ont pas été respectées.

Vous ne pouvez télécharger, exporter ou réexporter ces informations qu'en total accord avec toutes les lois et règlements applicables dans votre pays, y compris les lois et règlements américains relatifs à l'exportation.

IBM N'OCTROIE AUCUNE GARANTIE SUR LE CONTENU DE CES PUBLICATIONS. LES PUBLICATIONS SONT LIVREES "EN L'ETAT" SANS AUCUNE GARANTIE EXPLICITE OU IMPLICITE. IBM DECLINE NOTAMMENT TOUTE RESPONSABILITE RELATIVE A CES PUBLICATIONS EN CAS DE CONTREFAÇON AINSI QU'EN CAS DE DEFAUT D'APTITUDE A L'EXECUTION D'UN TRAVAIL DONNE.

Index

Α

ajout de connexions IBM SPSS Modeler Server 9 ajuster les flux à la vue 19 analyse de vente au détail 223 analyse discriminante carte territoriale 242 lambda de Wilk 240 matrice de structure 241 méthode détaillée étape par étape 239 table de classification 243 valeurs propres 240 analyse du panier d'achats 329 annuler 16 arrêter l'exécution 16

В

barre d'outils 16 bouton central de la souris simulation 20

С

canevas 12 carte territoriale analyse discriminante 242 champs classement par importance 93 filtrage 93 sélection pour analyse 93 classement des prédicteurs 93 classes 16 CLEM introduction 22 codages de variables catégorielles dans la régression de Cox 303 coller 16 connexion à IBM SPSS Modeler Server 8 connexion unique 8 connexions à IBM SPSS Analytic Server 10 à IBM SPSS Modeler Server 8,9 cluster de serveurs 9 coordinateur de processus 9 COP 9 copier 16 couper 16 courbes de risque dans la régression de Cox 307 courbes de survie dans la régression de Cox 307 CRISP-DM 16

D

documentation 3

données

affichage 80 lecture 77 manipulation 85 modélisation 88, 90, 91 données de survie avec censure par intervalle dans les modèles linéaires généralisés 245 données de survie groupées dans les modèles linéaires généralisés 245 Down Search modèles Liste de décision 110

Ε

estimations de paramètres dans les modèles linéaires généralisés 252, 264, 278, 287 Excel connexion aux modèles Liste de décision 123 modification des modèles Liste de décisions 129 exemples analyse de vente au détail 223 analyse discriminante 233 analyse du panier d'achats 329 classification d'échantillon de cellules 289 estimation d'une offre de nouveau véhicule 337 Guide des applications 3 KNN 337 noeud Recoder 101 réduction de la longueur des chaînes 101 réduction de la longueur des chaînes d'entrée 101 régression logistique multinomiale 133, 141 Réseau Bayésien 205, 213 surveillance d'état 227 SVM 289 télécommunications 133, 141, 153, 171, 233 ventes sur catalogue 179 vue d'ensemble 4 exemples d'application 3

F

fenêtre principale 12 filtrage 88 filtrage des prédicteurs 93 flux 7, 12 ajuster à la vue 19 création 77

G

générateur de formules 85 génération de scripts 22 gestionnaires 15

IBM SPSS Analytic Server connexion 10 connexions multiples 10 IBM SPSS Modeler 1, 12 démarrage 7 documentation 3 exécution depuis la ligne de commande 7 vue d'ensemble 7 IBM SPSS Modeler Server 1 ID utilisateur 8 mot de passe 8 nom d'hôte 8, 9 nom de domaine (Windows) 8 numéro de port 8,9 icônes définition des options 19 ID utilisateur IBM SPSS Modeler Server 8 importance classement des prédicteurs 93 Impression 21 flux 19 introduction IBM SPSS Modeler 7

lambda de Wilk analyse discriminante 240 ligne de commande démarrage d'IBM SPSS Modeler 7

Μ

matrice de structure analyse discriminante 241 méthode détaillée étape par étape analyse discriminante 239 dans la régression de Cox 304 Microsoft Excel connexion aux modèles Liste de décision 123 modification des modèles Liste de décisions 129 modèles de causalité temporelle étude de cas 347 tutoriel 347 modèles linéaires généralisés estimations de paramètres 252, 264, 278, 287 procédures apparentées 270, 281, 288 modèles linéaires généralisés (suite) qualité d'ajustement 276, 280 régression de Poisson 271 test composite 277 tests des effets du modèle 250, 262, 277 modèles Liste de décision connexion à Excel 123 enregistrement des informations de session 131 exemple d'application 107 génération 131 mesures personnalisées avec Excel 123 modification du modèle Excel 129 modèles Sélection de fonction 93 modélisation 88, 90, 91 mot de passe IBM SPSS Analytic Server 10 IBM SPSS Modeler Server 8 moyennes des covariables dans la régression de Cox 306

Ν

noeud Analyse 91 noeud dériver 85 noeud Liste de décision exemple d'application 107 noeud Modèle de réponse en auto-apprentissage création du flux 194 exemple d'application 193 exemple de création de flux 194 navigation dans le modèle 198 noeud MRAA création du flux 194 exemple d'application 193 exemple de création de flux 194 navigation dans le modèle 198 noeud Relations 83 Noeud Sélection de fonction classement des prédicteurs 93 filtrage des prédicteurs 93 importance 93 noeud Table 80 noeuds 7 noeuds Graphiques 83 noeuds source 77 nom d'hôte IBM SPSS Modeler Server 8, 9 nom de domaine (Windows) IBM SPSS Modeler Server 8 nuggets définis 15 numéro de port IBM SPSS Modeler Server 8, 9

0

observations censurées dans la régression de Cox 302

Ρ

palette de modèles générés 15 palettes 12 prédicteurs classement par importance 93 filtrage 93 sélection pour analyse 93 préparation 85 programmation visuelle 12 projets 16

Q

qualité d'ajustement dans les modèles linéaires généralisés 276, 280

R

raccourcis clavier 20 recherche à faible probabilité modèles Liste de décision 110 recherche de connexions dans COP 9 redimensionnement 18 réduction 18 régression binomiale négative dans les modèles linéaires généralisés 279 Régression de Cox codages de variables catégorielles 303 courbe de risque 307 courbe de survie 307 observations censurées 302 sélection des variables 304 régression de Poisson dans les modèles linéaires généralisés 271 régression gamma dans les modèles linéaires généralisés 283 répertoire temporaire 11 reste modèles Liste de décision 110

S

segments exclusion de l'évaluation 110 modèles Liste de décision 110 serveur ajout de connexions 9 connexion 8 recherche de serveurs dans COP 9 sessions IBM SPSS Modeler multiples 11 sortie 15 souris utilisation dans IBM SPSS Modeler 20 surveillance d'état 227

Т

table de classification analyse discriminante 243 tâches d'exploration modèles Liste de décision 110 test composite dans les modèles linéaires généralisés 277 tests composites dans la régression de Cox 304 tests des effets du modèle dans les modèles linéaires généralisés 250, 262, 277 titulaire IBM SPSS Analytic Server 10 touches de raccourci 20

U

URL IBM SPSS Analytic Server 10

V

valeurs propres analyse discriminante 240 var. noeud de fichiers 77 visualiseur Liste de décision 110 visualiseur Liste interactive exemple d'application 110 Panneau d'aperçu 110 utilisation 110

Ζ

zoom 16



Imprimé en France