
Qualité d'ajustement d'arbres d'induction

Gilbert Ritschard* — Djamel A. Zighed**

* Département d'économétrie, Université de Genève
bd du Pont-d'Arve 40, CH-1211 Genève 4
gilbert.ritschard@themes.unige.ch

** Laboratoire ERIC, Université Lyon 2
Bat. L, C.P. 11, F-69676 Bron Cédex
zighed@univ-lyon2.fr

RÉSUMÉ. Cet article discute des possibilités de mesurer la qualité de l'ajustement d'arbres d'induction aux données comme cela se fait classiquement pour les modèles statistiques. On montre comment adapter aux arbres d'induction les statistiques du khi-2, notamment celle du rapport de vraisemblance utilisée dans le cadre de la modélisation de tables de contingence. Cette statistique permet de tester l'ajustement du modèle, mais aussi l'amélioration de l'ajustement qu'apporte la complexification de l'arbre. On en déduit également des formes adaptées des critères d'information AIC et BIC qui permettent de sélectionner le meilleur arbre en terme de compromis entre ajustement et complexité.

ABSTRACT. This paper is concerned with the fit of induction trees. Namely, we explore the possibility to measure the goodness-of-fit as it is classically done in statistical modeling. We show how Chi-square statistics and especially the Log-likelihood Ratio statistic that is abundantly used in the modeling of cross tables, can be adapted for induction trees. Not only is the Log-likelihood Ratio statistic suited for testing the fit. It allows also to test the significance of the fit improvement provided by the complexification of a tree. In addition, we derive from it adapted forms of the Akaike (AIC) and Bayesian (BIC) information criteria that prove useful in selecting the best compromise tree between fit and complexity.

MOTS-CLÉS : arbre d'induction, qualité d'ajustement, tests du khi-2, comparaison de modèles

KEYWORDS: Induction tree, goodness-of-fit, Chi-square tests, models comparison

1. Introduction

Les arbres d'induction [BRE 84][QUI 93][ZIG 00] sont l'un des outils les plus populaires d'apprentissage supervisé. Ils consistent à rechercher par éclatements successifs de sommets, une partition de l'ensemble des combinaisons de valeurs des prédicteurs optimale pour prédire la variable réponse. La prédiction se fait simplement en choisissant, dans chaque classe de la partition obtenue, la modalité la plus fréquente de la variable à prédire. Bien que leur utilisation première soit la classification, les arbres d'induction fournissent une description de la façon dont la distribution de la variable à prédire est conditionnée par les valeurs des prédicteurs. Ils nous indiquent par exemple comment la répartition entre clients solvables et insolvable est influencée par les attributs âge, sexe, niveau d'éducation, profession, etc. En ce sens les arbres d'induction sont donc des outils de modélisation de l'influence des prédicteurs sur la variable à prédire au même titre que par exemple la régression linéaire, la régression logistique ou la modélisation log-linéaire de tables de contingence multi-dimensionnelles. C'est essentiellement à cet aspect de modélisation descriptive, et en particulier à l'évaluation de la qualité de la description fournie par un arbre induit que nous nous intéressons dans cet article.

En modélisation statistique, qu'il s'agisse de régression linéaire, d'analyse discriminante, de régression logistique ou plus généralement de modèle linéaire généralisé (GLM), il est d'usage d'évaluer la qualité d'ajustement du modèle, c'est-à-dire la qualité de la description fournie par le modèle, avec des mesures descriptives telles que le coefficient de détermination R^2 ou des pseudo R^2 , et avec des statistiques de test telles que les khi-2 du score test, de Wald ou du rapport de vraisemblance. Parmi ces dernières, la statistique du rapport de vraisemblance jouit d'une propriété d'additivité qui permet d'évaluer également la pertinence statistique de la simplification d'un modèle de référence par renforcement de contraintes sur ses paramètres, ou, si l'on regarde les choses dans l'autre sens, la significativité statistique de la complexification résultant de l'ajout de paramètres à un modèle donné.

Le cas particulier des tests "omnibus", où l'on teste globalement l'apport des facteurs explicatifs par rapport à un modèle de référence simple — le modèle avec la constante comme seule variable explicative dans le cas de la régression, le modèle d'indépendance dans le cas des modèles log-linéaires — retiendra notre attention. Dans le cas des arbres d'induction, le modèle de référence est naturellement l'arbre de niveau 0 constitué par le seul nœud initial.

Une des difficultés principales à laquelle on se heurte dans la pratique des arbres ou graphes d'induction est le fort degré de complexité des arbres mis en évidence. Il nous paraît alors souhaitable de pouvoir disposer aussi de critères tels que le critère d'information d'Akaike (AIC) [AKA 73] ou le critère d'information bayésien (BIC) de Schwarz [SCH 78] [KAS 95]. Ces critères sont une combinaison de la qualité d'ajustement (statistique du rapport de vraisemblance) et d'une mesure de la complexité. Ils s'avèrent ainsi une aide précieuse pour arbitrer entre complexité et qualité d'ajustement dans la sélection de modèles.

L'article est organisé comme suit. La section 2 situe la place des mesures de qualité d'ajustement parmi les mesures classiques de qualité d'un arbre d'induction. La section 3 précise le concept d'ajustement considéré. La section 4 est consacrée aux critères de qualité d'ajustement. On montre comment les statistiques du khi-2 de Pearson et du rapport de vraisemblance utilisées dans le cadre de la modélisation de tables de contingence peuvent s'adapter aux arbres d'induction. Nous discutons ensuite l'amélioration qu'apporte un modèle par rapport à un modèle de référence. On montre comment exploiter la différence des statistiques G^2 du rapport de vraisemblance pour tester la significativité statistique du gain d'information et proposons des pseudo R^2 . Nous montrons également comment appliquer les critères d'information d'Akaike (AIC) et bayésien (BIC) aux arbres d'induction et discutons leur intérêt pour guider le choix entre modèles de complexité variable. La section 5 illustre l'utilisation des critères proposés dans un cas concret. Finalement, la conclusion fait l'objet de la section 6 où nous mentionnons des pistes de développements futurs de la démarche initiée dans cet article.

2. Arbre d'induction et mesures classiques de qualité

Avant de se concentrer sur la qualité d'ajustement, nous rappelons brièvement le principe des arbres d'induction et leurs critères usuels de qualité. Ceci dans le but de pouvoir mieux situer le rôle des mesures de qualité d'ajustement. Pour une discussion plus détaillée le lecteur peut par exemple se référer à [ZIG 00].

2.1. Rappel du principe des arbres d'inductions

L'objectif est de construire une règle qui permette, à partir de la connaissance d'un vecteur d'attribut $\mathbf{x} = (x_1, \dots, x_p)$, de prédire une variable réponse y , ou si l'on préfère de classer les cas selon les états de la variable y . La construction de la règle se fait en deux temps. Dans un premier temps, on détermine une partition des valeurs possibles de \mathbf{x} telle que la distribution de la réponse y soit la plus pure possible dans

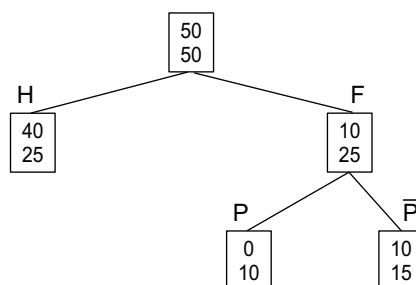


Figure 1. Arbre induit

chaque classe de la partition, ou de façon plus ou moins équivalente la plus différente possible d'une classe à l'autre. La règle consiste ensuite à attribuer à chaque cas la valeur de y la plus fréquente dans sa classe.

Les arbres d'induction déterminent la partition par éclatements successifs des sommets. En partant du sommet initial, ils recherchent l'attribut qui permet le meilleur éclatement selon un critère qui peut être par exemple le gain d'entropie (C4.5, Sipina) ou la significativité d'un χ^2 (CHAID). L'opération est répétée à chaque nouveau sommet jusqu'à ce qu'un critère d'arrêt, une taille minimale du sommet par exemple, soit atteint. Le résultat est un arbre tel que celui présenté à la figure 1.

2.2. Taux d'erreur

Le taux d'erreur de classification, c'est-à-dire le pourcentage de cas mal classés est peut-être le critère de qualité le plus utilisé. S'agissant de classification, c'est évidemment la performance en généralisation qui importe, c'est-à-dire la performance pour des cas n'ayant pas servi à l'apprentissage. C'est pourquoi il convient de calculer le taux d'erreur sur un échantillon de validation différent de l'échantillon d'apprentissage.

La décomposition biais/variance/erreur résiduelle du taux d'erreur [GEU 02] fournit des indications précieuses sur la part de l'erreur due à la variabilité de l'échantillon d'apprentissage. Ces décompositions ne sont cependant quantifiables que numériquement sur la base de simulations ou de méthodes de rééchantillonnages ce qui limite leur utilisation systématique.

2.3. Qualité des partitions

La qualité d'une partition est d'autant meilleure que les sommets terminaux sont purs, c'est-à-dire qu'ils ont des distributions le plus proche possible de la distribution dégénérée qui donne un poids de un à l'une des modalités et zéro aux autres.

Les mesures de qualité des partitions obtenues sont notamment discutées dans [ZIG 00], chap. 9. Parmi ces mesures on peut mentionner :

- Les mesures issues de la *théorie de l'information* et qui consistent essentiellement à mesurer l'entropie pour la partition considérée.
- Celles qui se fondent sur des *distances entre distributions de probabilités*. Le principe consiste ici à vérifier que les distributions diffèrent le plus possible d'une classe à l'autre. La démarche est simple dans le cas de deux classes. Pour le cas plus général d'un nombre quelconque de classes, les solutions proposées restent de portée assez limitée.
- Enfin, des mesures qui s'appuient sur des *indices d'association*. L'idée est ici que la partition est d'autant meilleure que l'association entre la classe et la variable

dépendante est forte. [RAK 98] proposent d'utiliser le degré de signification du τ de Goodman-Kruskal pour prendre en compte également la taille des échantillons.

2.4. Complexité

Un des intérêts majeurs souvent avancé des arbres et graphes d'induction est la facilité de leur interprétation. Ceci est vrai tant que la complexité de l'arbre reste limitée, d'où l'intérêt de mesurer cette complexité. Les indicateurs couramment utilisés sont :

- *Le nombre de sommets terminaux.* Ce critère correspond au nombre de règles de prédiction ainsi qu'au nombre de classes de la partition finale.
- *Le nombre de nœuds.* Ce critère est évidemment lié à la procédure de construction. Il sera en général plus élevé pour un arbre binaire qui tend à multiplier les sommets intermédiaires. Ce critère reflète bien la complexité visuelle de l'arbre induit.
- *La profondeur de l'arbre.* Elle dépend également de la procédure de construction et reflète la complexité visuelle de l'arbre.
- *La longueur des messages.* Traduit la complexité des règles d'affectation aux différentes classes.

A la section 4.3, nous verrons que le nombre de paramètres généralement utilisé comme mesure de complexité en modélisation statistique est également un concept pertinent pour les arbres d'induction.

2.5. Place des mesures de qualité d'ajustement

Les critères d'évaluation de la qualité d'ajustement discutés dans les sections suivantes concernent la capacité d'un arbre à reproduire la distribution de la réponse y pour les individus ayant un profil x donné. En d'autres termes, on s'intéresse à la qualité de reproduction de la table de contingence qui croise la variable réponse y avec l'ensemble des prédicteurs. Il s'agit donc de la qualité descriptive de l'arbre par opposition à sa performance en classification.

Par rapport à la typologie précédente, les mesures considérées ci-après s'inscrivent essentiellement dans l'optique qualité des partitions. En tant qu'instruments d'arbitrage entre ajustement et complexité, les critères d'information AIC et BIC concernent cependant également la complexité .

La mesure de la qualité d'ajustement donne des indications complémentaires aux critères usuels mentionnés ci-dessus. Essentiellement, il s'agit de rendre compte de deux aspects :

- 1) l'aptitude de l'arbre induit à décrire la distribution de la variable réponse conditionnellement aux valeurs prises par les prédicteurs,

2) le gain d'information qu'apporte l'arbre induit par rapport au nœud initial où l'on ne tient pas compte des prédicteurs.

3. Concept de qualité d'ajustement d'un arbre

3.1. Table cible et table prédite

De façon générale, la qualité d'ajustement d'un modèle se réfère à sa capacité à reproduire les données. Dans le cas de la prédiction quantitative d'une variable Y , par exemple dans le cas de la régression linéaire, l'objectif est clair. Il s'agit d'obtenir des valeurs prédites \hat{y}_α qui s'ajustent le mieux possible aux valeurs observées y_α , pour $\alpha = 1, \dots, n$, n étant le nombre d'observations. De même, dans l'optique de la classification, les états prédits \hat{y}_α doivent correspondre le plus souvent possible aux vraies valeurs y_α . Le taux d'erreur est dans ce cas un indicateur naturel de qualité d'ajustement.

Dans certaines situations, en particulier en sciences sociales, les arbres ou graphes d'induction sont utilisés plus dans une optique descriptive que prédictive comme outil de mise en évidence des principaux déterminants de la variable à prédire. Ils sont utilisés comme outil d'aide à la compréhension de phénomènes et non pas comme outil de classification.

Ce ne sont plus alors les états particuliers y_α que l'on cherche à reproduire. Pour comprendre comment les prédicteurs interagissent sur la variable réponse Y , il convient en effet d'examiner comment la distribution de Y change avec le profil \mathbf{x} . Dans cette optique, la qualité d'ajustement considérée ici, se réfère à la qualité de la reproduction de l'ensemble des distributions conditionnelles empiriques.

En admettant que les prédicteurs prennent un nombre fini de valeurs, l'ensemble des distributions conditionnelles est représentable sous forme d'une table de contingence \mathbf{T} croisant y avec une variable composite définie par le croisement de tous les prédicteurs. Le tableau 1 est un exemple d'une telle table dans le cas où la variable à prédire est le statut marital et les prédicteurs le genre et le secteur d'activité. Le nombre de lignes de \mathbf{T} est le nombre ℓ d'états de la variable Y . De même, si chaque prédicteur x_j , $j = 1, \dots, p$ a c_j valeurs différentes, le nombre c de colonnes de \mathbf{T} est au plus le produit des c_j , soit : $c \leq \prod_j c_j$, certaines combinaisons de valeurs des attributs pouvant être structurellement impossibles. On s'intéresse donc à la capacité de l'arbre à reproduire cette table \mathbf{T} .¹

Un élément de la table \mathbf{T} est noté n_{ij} et représente le nombre de cas avec profil \mathbf{x}_j qui dans les données prennent la valeur y_i de la variable réponse. On note $\hat{\mathbf{T}}$ la table prédite à partir d'un arbre et \hat{n}_{ij} désigne un élément générique de cette table.

1. Il n'est pas sans intérêt de relever comme nous l'avons fait dans [RIT 03], que dans cette optique de reconstitution de T , les arbres peuvent constituer un outil alternatif et complémentaire à la modélisation log-linéaire des tables de contingence multidimensionnelle.

Tableau 1. Exemple de table de contingence \mathbf{T}

marié	homme			femme			total
	primaire	secondaire	tertiaire	primaire	secondaire	tertiaire	
non	11	14	15	0	5	5	50
oui	8	8	9	10	7	8	50
total	19	22	24	10	12	13	100

Formellement, la qualité d'ajustement se réfère ainsi à la divergence entre les tables \mathbf{T} et $\hat{\mathbf{T}}$.

Il reste à préciser comment l'on déduit la table estimée $\hat{\mathbf{T}}$ à partir d'un arbre. On utilise pour cela le modèle de reconstruction suivant où l'on note \mathbf{T}_j la j -ème colonne de \mathbf{T} :

$$\hat{\mathbf{T}}_j = n a_j \hat{\mathbf{p}}_{|j}, \quad j = 1, \dots, c \quad (1)$$

Les paramètres sont le nombre total n de cas, les proportions a_j de cas par colonne $j = 1, \dots, c$, et les c vecteurs de probabilités $\mathbf{p}_{|j} = \mathbf{p}(Y|j)$ correspondant à la distribution de Y dans chaque colonne j de la table. Nous verrons que chaque arbre donne lieu à des estimations différentes des vecteurs $\mathbf{p}_{|j}$ et par suite à une estimation $\hat{\mathbf{T}}$ différente.

3.2. Arbre saturé et arbre étendu

En modélisation statistique, on appelle modèle saturé un modèle avec le nombre maximal de paramètres libres qui peuvent être estimés à partir des données. En modélisation log-linéaire de tables de contingence multi-dimensionnelles, le modèle saturé permet de reproduire exactement la table modélisée. Par analogie, on introduit le concept d'arbre saturé qui permet de reproduire exactement la table \mathbf{T} .

Définition 1 (Arbre saturé) *Pour des variables prédictives catégorielles, on appelle arbre saturé, un arbre qui résulte de tous les éclatements successifs possibles selon les modalités des variables prédictives.*

Les distributions conditionnelles aux feuilles de l'arbre saturé sont estimées par les vecteurs de fréquences relatives $\mathbf{p}_{|j}$, soit les vecteurs d'éléments $n_{ij}/n_{.j}$, $i = 1, \dots, \ell$, le $n_{.j}$ étant le total de la j -ème colonne de \mathbf{T} .

L'arbre saturé n'est pas unique, des variantes étant possibles selon l'ordre dans lequel les variables sont prises en compte. Tous les arbres saturés conduisent cependant aux mêmes feuilles (sommets terminaux). Ces feuilles correspondent aux colonnes de la table de contingence \mathbf{T} .

Par exemple, la figure 2 donne l'arbre saturé correspondant au cas du tableau 1 où l'on cherche à prédire le statut marital connaissant le genre ($H =$ homme, $F =$ femme)

et le secteur d'activité (P = primaire, S = secondaire, T = tertiaire). Les distributions conditionnelles sont :

$$\mathbf{P}_{|HP} = \begin{pmatrix} 11/19 \\ 8/19 \end{pmatrix}, \quad \mathbf{P}_{|HS} = \begin{pmatrix} 14/22 \\ 8/22 \end{pmatrix}, \quad \mathbf{P}_{|HT} = \begin{pmatrix} 15/24 \\ 9/24 \end{pmatrix},$$

$$\mathbf{P}_{|FP} = \begin{pmatrix} 0/10 \\ 10/10 \end{pmatrix}, \quad \mathbf{P}_{|FS} = \begin{pmatrix} 5/12 \\ 7/12 \end{pmatrix}, \quad \mathbf{P}_{|FT} = \begin{pmatrix} 5/13 \\ 8/13 \end{pmatrix}$$

Notons qu'un algorithme de génération d'arbre ne peut en général générer un arbre maximal théorique que si i) toutes les cellules de la table de contingence des variables prédictives sont non vides et si ii) la distribution de la variable dépendante est différente dans chacune des cellules.

Pour comparer les distributions des feuilles de l'arbre induit à celles de l'arbre maximal, on doit étendre l'arbre induit pour obtenir des feuilles de même définition et en particulier en même nombre que l'arbre saturé.

Définition 2 (Extension maximale d'un arbre d'induction) *Pour des variables prédictives catégorielles, on appelle extension maximale de l'arbre induit ou arbre induit étendu, l'arbre qui résulte de tous les éclatements successifs possibles de ses sommets terminaux. On applique aux feuilles de l'extension la distribution $\mathbf{p}_{|k}^a$ du nœud terminal parent de l'arbre initial. On note $\hat{\mathbf{p}}_{|i}$ les distributions conditionnelles des feuilles de l'extension maximale de l'arbre.*

Par exemple, si l'arbre d'induction obtenu est l'arbre avec les sommets blancs de la figure 3, son extension maximale s'obtient en ajoutant les sommets gris et en répartissant dans ceux-ci l'effectif selon la distribution du sommet dont ils sont issus.

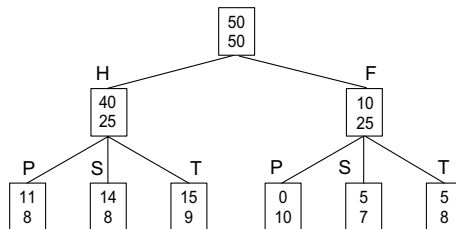


Figure 2. Arbre saturé

Les distributions des six sommets terminaux de l'extension se déduisent de celles des trois sommets terminaux de l'arbre induit, soit pour notre exemple :

$$\hat{\mathbf{P}}_{|HP} = \hat{\mathbf{P}}_{|HS} = \hat{\mathbf{P}}_{|HT} = \mathbf{P}_{|H}^a = \begin{pmatrix} 40/65 \\ 25/65 \end{pmatrix}$$

$$\hat{\mathbf{P}}_{|FP} = \mathbf{P}_{|FP}^a = \begin{pmatrix} 0/10 \\ 10/10 \end{pmatrix}$$

$$\hat{\mathbf{P}}_{|FS} = \hat{\mathbf{P}}_{|FT} = \mathbf{P}_{|F\bar{P}}^a = \begin{pmatrix} 10/25 \\ 15/25 \end{pmatrix}$$

Les feuilles terminales de l'extension de l'arbre induit nous donnent la prédiction $\hat{\mathbf{T}}$ de la table \mathbf{T} , la table 2

Tableau 2. Exemple de table de contingence prédite $\hat{\mathbf{T}}$

marié	homme			femme			total
	primaire	secondaire	tertiaire	primaire	secondaire	tertiaire	
non	11.7	13.5	14.8	0	4.8	5.2	50
oui	7.3	8.5	9.2	10	7.2	7.8	50
total	19	22	24	10	12	13	100

4. Qualité d'ajustement d'un arbre induit

Nous montrons comment adapter aux arbres d'induction les mesures habituellement utilisées en modélisation statistique. Nous commençons avec les tests d'ajustement du khi-2

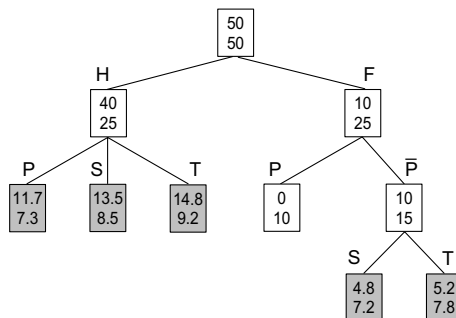
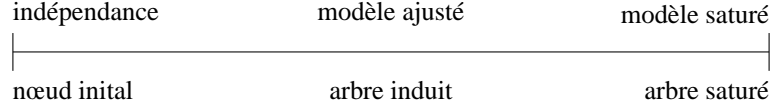


Figure 3. Arbre induit (sommets blancs) et son extension maximale

Comme dans la modélisation multinomiale log-linéaire [AGR 90], les statistiques de test du khi-2 sur la divergence entre $\hat{\mathbf{T}}$ et \mathbf{T} permettent de juger de l'écart entre l'arbre induit (modèle ajusté) et l'arbre saturé (modèle saturé).



Plutôt que de chercher à savoir à quelle distance l'on se trouve du modèle saturé correspondant à la partition la plus fine, il peut être utile de savoir ce que l'on a gagné par rapport à la situation d'indépendance où l'on ne tient pas compte des prédicteurs. Cet aspect correspond à l'écart entre l'arbre induit et le modèle d'indépendance représenté par l'arbre constitué du seul nœud initial. Dans cette optique, nous introduisons les tests sur l'amélioration de l'ajustement et les pseudo R^2 . Enfin, nous établissons des formes des critères d'information AIC et BIC adaptées aux arbres d'induction.

4.1. Statistiques du khi-2 pour arbres d'induction

L'objectif est de mesurer la divergence entre \mathbf{T} et $\hat{\mathbf{T}}$ avec une statistique permettant d'évaluer la significativité de l'écart observé.

Les $\mathbf{p}_{|j}$ étant estimés par le maximum de vraisemblance, on peut simplement appliquer les statistiques X^2 du khi-2 de Pearson et G^2 du maximum de vraisemblance :

$$X^2 = \sum_{i=1}^{\ell} \sum_{j=1}^c \frac{(n_{ij} - \hat{n}_{ij})^2}{\hat{n}_{ij}} \quad (2)$$

$$G^2 = 2 \sum_{i=1}^{\ell} \sum_{j=1}^c n_{ij} \ln \left(\frac{n_{ij}}{\hat{n}_{ij}} \right) \quad (3)$$

Sous l'hypothèse que le modèle est correct et sous certaines conditions de régularité, voir par exemple [BIS 75] chap. 14, ces statistiques suivent une même distribution du χ^2 avec pour degrés de liberté le nombre de cellules de la table de contingence moins le nombre de paramètres linéairement indépendants du modèle (1) de prédiction de \mathbf{T} . Précisons alors le nombre de paramètres linéairement indépendants.

Un arbre induit non saturé définit une partition de l'ensemble \mathcal{X} des profils \mathbf{x} possibles. Chacun de ses q sommets terminaux correspond donc à un sous-ensemble $\mathcal{X}_k \subseteq \mathcal{X}$, $k = 1, \dots, q$ de profils x_j pour lequel on impose la contrainte

$$\hat{\mathbf{p}}_{|j} = \mathbf{p}_{|k}^a \quad \text{pour tout } \mathbf{x}_j \in \mathcal{X}_k \quad k = 1, \dots, q \quad (4)$$

où $\mathbf{p}_{|k}^a$ désigne la distribution dans le sommet terminal k de l'arbre induit.

Chaque vecteur $\mathbf{p}_{|k}^a$ contient $\ell - 1$ termes indépendants et il y a $q - 1$ distributions conditionnelles $\mathbf{p}_{|k}^a$ indépendantes. On en déduit le décompte ci-dessous des paramètres indépendants pour un nombre q de fixé de sommets terminaux.²

paramètres	nombre	dont indépendants
$p_{i j}, i = 1, \dots, \ell, j = 1, \dots, c$	$c\ell$	$q(\ell - 1)$
$a_j, j = 1, \dots, c$	c	$c - 1$
n	1	1
Total	$c\ell + \ell + c + 1$	$q\ell - q + c$

En retranchant au nombre $c\ell$ de cellules de \mathbf{T} le nombre $q(\ell - 1) + c$ de paramètres indépendants, on obtient les degrés de liberté d_M de l'arbre induit, soit

$$\text{degrés de liberté} = d_M = (c - q)(\ell - 1) .$$

Notons que ce nombre correspond au nombre de contraintes (4). Pour le modèle d'indépendance on a $q = 1$ et l'on retrouve la valeur usuelle des degrés de liberté du test d'indépendance, soit $d_I = (c - 1)(\ell - 1)$. De même, pour l'arbre saturé on a $q = c$ et donc $d_S = 0$.

Sous réserve des conditions de régularité, les statistiques X^2 et G^2 de tables associées à des arbres suivent donc, lorsque le modèle est correct, une distribution du χ^2 avec $(c - q)(\ell - 1)$ degrés de liberté.

Pour l'arbre de la figure 1, on trouve par exemple, $X^2 = 0.1823$ et $G^2 = 0.1836$. On a $c = 6$, $q = 3$ et $\ell = 2$, et donc $d_M = (6 - 3)(2 - 1) = 3$ degrés de liberté. Les valeurs des statistiques sont très petites et, avec un degré de signification de l'ordre de 98% dans les deux cas, indiquent clairement que $\hat{\mathbf{T}}$ ajuste de façon satisfaisante \mathbf{T} . La qualité de l'ajustement de l'arbre induit aux données est donc dans ce cas excellente.

Théoriquement, les statistiques X^2 de Pearson (2) et G^2 du rapport de vraisemblance (3) devraient permettre de tester si l'arbre induit s'ajuste de façon satisfaisante aux données. Il est bien connu cependant que la portée du test reste limitée lorsque l'échantillon est grand, le moindre écart devenant alors statistiquement significatif. Par ailleurs, les conditions de régularité requises pour que les statistiques soient distribuées selon une loi du χ^2 , et en particulier les conditions d'intériorité (pas d'effectifs attendus nuls), sont difficilement tenables lorsque l'arbre saturé compte un grand nombre de feuilles.

Plus intéressante nous semble être l'utilisation de la statistique G^2 pour comparer des modèles imbriqués, un modèle M_2 (restreint) étant inclus dans M_1 (non restreint) si l'espace de ses paramètres est un sous-ensemble de M_1 , c'est-à-dire, en d'autres

2. Bien que q soit induit des données, on raisonne ici conditionnellement à q comme cela se fait en modélisation log-linéaire où les interactions prises en compte sont également induites des données par le biais du processus de sélection du modèle.

termes, si les paramètres de M_2 s'obtiennent en imposant des contraintes sur ceux du modèle M_1 . En effet dans ce cas la déviance entre les deux modèles est (voir par exemple [AGR 90], p.211 ou [POW 00], p. 105) :

$$G^2(M_2|M_1) = G^2(M_2) - G^2(M_1) \quad (5)$$

qui, sous l'hypothèse que M_2 est correct lorsque M_1 l'est, est approximativement distribuée selon une loi du χ^2 avec pour degrés de liberté la différence $d_2 - d_1$ des degrés de liberté des modèles M_2 et M_1 .

Cette dernière propriété permet en particulier de tester la significativité d'un éclatement. La déviance entre le modèle après l'éclatement et celui avant l'éclatement nous renseigne en effet sur la pertinence statistique de cet éclatement. Par exemple, si M_1 est l'arbre induit de la figure 1 et M_2 l'arbre avant l'éclatement du sommet « femme ». On a $G^2(M_1) = 0.18$ avec 3 degrés de liberté et $G^2(M_2) = 8.41$ avec 4 degrés de liberté. La déviance est alors

$$G^2(M_2|M_1) = 8.41 - 0.18 = 8.23 \quad \text{avec} \quad d_2 - d_1 = 4 - 3 = 1$$

Son degré de signification (p -valeur) est 0.4%, donc inférieur au seuil généralement admis de 5%, ce qui indique que l'éclatement est statistiquement significatif.

4.2. Comparaison avec un modèle de référence

Cette section est consacrée aux indicateurs de type R^2 qui mesurent le gain relatif de qualité d'un arbre induit M par rapport à l'arbre trivial constitué par le seul nœud initial. On discute successivement le pourcentage d'amélioration du taux d'erreur, le pourcentage de réduction de l'entropie, l'amélioration de l'ajustement et les pseudo R^2 . Les mesures considérées ici pour la comparaison entre le graphe induit et le nœud initial se généralisent aisément, bien que nous ne les traitons pas explicitement, au cas général de la comparaison de deux modèles dont l'un est inclus dans l'autre.

4.2.1. Remarque sur le pourcentage de réduction du taux d'erreur

Bien qu'on ne s'intéresse pas ici à la prédiction de valeurs individuelles, nous aimerions souligner que l'idée de comparer la performance du modèle avec le modèle qui ne tient pas compte des prédicteurs est également pertinente en terme de taux d'erreur. On peut noter que, sur l'échantillon d'apprentissage, le pourcentage de réduction de l'erreur de classification correspond à la mesure d'association $\lambda_{y|\text{partition}}$ de Goodman-Kruskal [GOO 54] entre la variable dépendante y et la partition définie par le graphe induit. Il existe une forme analytique de la variance asymptotique de cette mesure qui permet d'en tester la significativité sous certaines conditions de régularité [GOO 72] [OLS 95]. Nous n'approfondissons pas cet aspect ici, notre objectif étant les qualités descriptives du graphe induit plutôt que sur ses qualités prédictives.

4.2.2. Gain d'information

Le gain d'information peut être mesuré par la réduction de l'entropie de la distribution de la variable dépendante que permet la connaissance des classes de la partition définie par l'arbre induit. Précisons que nous nous intéressons à la réduction globale que permet l'arbre par rapport à la distribution marginale, c'est-à-dire celle dans le nœud initial. Le gain discuté ici se distingue donc du gain partiel et conditionnel à un nœud que certains algorithmes cherchent à maximiser à chaque étape de construction de l'arbre.

Le gain d'information relativement au nœud initial est mesuré par exemple par les deux indicateurs suivants :

$$\hat{\tau}_{M|I} = \frac{n \sum_i \sum_j \frac{\hat{n}_{ij}^2}{n_j} - \sum_i n_i^2}{n^2 - \sum_i n_i^2} \quad (6)$$

$$\hat{u}_{M|I} = \frac{\sum_i \frac{n_{i.}}{n} \log_2 \frac{n_{i.}}{n} - \sum_j \frac{n_{.j}}{n} \sum_i \frac{\hat{n}_{ij}}{n_j} \log_2 \frac{\hat{n}_{ij}}{n_j}}{\sum_i \frac{n_{i.}}{n} \log_2 \frac{n_{i.}}{n}} \quad (7)$$

Le $\hat{\tau}_{M|I}$ est la deuxième mesure d'association nominale proposée par Goodman et Kruskal [GOO 54]. Il mesure la proportion de réduction de l'entropie quadratique $H_Q(\mathbf{p}) = \sum_i p_i(1 - p_i)$, connue aussi comme l'indice de variation de Gini. Le $\hat{u}_{M|I}$ est connu en statistique sous le nom de coefficient d'incertitude de Theil [THE 67] [THE 70]. Il mesure la proportion de réduction de l'entropie de Shannon $H_S(\mathbf{p}) = -\sum_i p_i \log_2 p_i$.

Pour l'arbre M de la figure 1, on trouve par exemple $\hat{\tau}_{M|I} = 0.145$ et $\hat{u}_{M|I} = 0.132$, valeurs qui indiquent une réduction d'entropie d'environ 14%. Si l'on compare ces valeurs à celles du modèle saturé, soit respectivement $\hat{\tau}_{S|I} = 0.146$ et $\hat{u}_{S|I} = 0.134$, il apparaît que l'arbre capte à peu près toute l'information que l'on peut tirer des attributs prédictifs retenus.

Pour tester la significativité statistique du gain d'information, soit les hypothèses $H_0 : \tau_{M|I} = 0$ et $H_0 : u_{M|I} = 0$, une possibilité est d'utiliser les variances asymptotiques que l'on trouve dans la littérature (voir par exemple [OLS 95]) et une approximation par la loi normale. Il est préférable cependant d'utiliser les transformations suivantes des indicateurs :

$$C(I|M) = (n-1)(\ell-1) \hat{\tau}_{M|I} \quad (8)$$

$$G^2(I|M) = \left(-2 \sum_i n_{i.} \log(n_{i.}/n) \right) \hat{u}_{M|I} \quad (9)$$

La première C est due à Light et Margolin [LIG 71] qui montrent que dans le cas où le tableau prédit est le modèle saturé ($M = S$), $C(I|S)$ est, sous l'hypothèse d'indépendance ($\tau_{S|I} = 0$), asymptotiquement distribuée comme un χ^2 à d_I degrés de liberté. Dans le cas d'un modèle M plus restrictif que S , il suffit d'adapter les

degrés de liberté. Ainsi, de façon générale $C(I|M)$ suit un χ^2 avec $d_I - d_M$ degrés de liberté.

La transformation du coefficient $\hat{u}_{M|I}$ montre que tester la significativité de $u_{M|I}$ est équivalent à tester la significativité de la différence d'ajustement entre I et M avec $G^2(I|M) = G^2(I) - G^2(M)$. Les deux transformations (8) et (9) suivent donc, sous H_0 , asymptotiquement la même loi du χ^2 à $d_I - d_M$ degrés de liberté. Le test avec ces statistiques est plus puissant qu'avec une approximation normale de $\hat{\tau}_{I|M}$ ou $\hat{u}_{I|M}$.

Pour notre exemple d'arbre induit, on trouve $C(I|M) = 14.32$ et $G^2(I|M) = 18.36$. Ces valeurs sont très grandes compte tenu des degrés de liberté $d_I - d_M = 5 - 3 = 2$. Elles confirment donc la significativité statistique du gain d'information de l'arbre par rapport au modèle d'indépendance.

4.2.3. Pseudo R^2

Dans une optique purement descriptive, on peut envisager, comme on le fait par exemple dans la modélisation log-linéaire, des pseudo R^2 qui mesurent la proportion de la déviance entre indépendance I et arbre saturé que l'arbre d'induction reproduit. On peut par exemple utiliser le pseudo R^2

$$R^2 = 1 - \frac{G^2(M)}{G^2(I)}$$

ou sa version corrigée des degrés de liberté

$$R_{\text{ajust}}^2 = 1 - \frac{G^2(M)/d_M}{G^2(I)/d_I}$$

Pour notre exemple, on a $G^2(I) = 18.55$, $d_I = 5$, $G^2(M) = .18$ et $d_M = 3$, d'où $R^2 = .99$ et $R_{\text{ajust}}^2 = .984$. Ces valeurs confirment que l'arbre capte presque le 100% de l'écart entre l'indépendance et la table cible représentée par l'arbre saturé.

Dans une optique de réduction d'entropie, nous proposons comme alternative au pseudo R^2 ci-dessus, de calculer la part de la proportion maximale de réduction d'entropie possible que l'on atteint avec l'arbre induit. La proportion maximale de réduction d'entropie est obtenue avec la partition la plus fine, c'est-à-dire le modèle saturé. Elle correspond à $\hat{\tau}_{S|I}$ pour l'entropie quadratique et $\hat{u}_{S|I}$ pour l'entropie de Shannon. Les parts de ces valeurs atteintes avec l'arbre induit sont donc :

$$R_{\tau}^2 = \frac{\hat{\tau}_{M|I}}{\hat{\tau}_{S|I}} \quad \text{et} \quad R_u^2 = \frac{\hat{u}_{M|I}}{\hat{u}_{S|I}}$$

Pour notre exemple, on obtient respectivement $R_{\tau}^2 = .993$ et $R_u^2 = .985$.

4.3. Ajustement et complexité

Dans le but de pouvoir arbitrer entre ajustement et complexité, on peut recourir aux critères d'information de Akaike (AIC) [AKA 73] ou au critère bayésien (BIC) de Schwarz [SCH 78] [KAS 95].

Dans notre cas, ces critères peuvent par exemple s'écrire :

$$\begin{aligned} \text{AIC}(M) &= G^2(M) + 2(q\ell - q + c) \\ \text{BIC}(M) &= G^2(M) + (q\ell - q + c) \log(n) \end{aligned}$$

Ici, la complexité est représentée par le nombre $q\ell - q + c$ de paramètres indépendants. Le critère BIC pénalise plus fortement la complexité que le critère AIC, la pénalisation augmentant avec le nombre de données n . Notons qu'il existe des formes alternatives du coefficient BIC. Raftery [RAF 95], par exemple, propose $BIC = G^2 - d \log(n)$, où d est le nombre de degrés de liberté, soit dans notre cas $d = (c - q)(\ell - 1)$. Comme ce nombre d diminue d'une unité chaque fois que l'on ajoute un paramètre indépendant, la pénalisation reste évidemment la même. Les deux formulations sont équivalentes à une translation $c\ell$ près.

Ces critères d'information offrent une alternative aux tests statistiques pour la sélection de modèles. Parmi plusieurs modèles, celui qui minimise le critère réalise le meilleur compromis entre ajustement et complexité. Le modèle qui minimise BIC en particulier, est, dans une approche bayésienne, optimal compte tenu de l'incertitude des modèles.

Pour illustrer l'utilisation de ces critères, on se propose de comparer notre arbre induit M avec la variante M^* où l'on éclate le sommet « femme » selon les trois secteurs P, S, I au lieu du partage binaire entre primaire P et non primaire \bar{P} . Dans les deux cas on a $n = 100$, $c = 6$ et $\ell = 2$. Pour M , on a $q = 3$ et donc $(q\ell - q + c) = 9$ et pour M^* , $q = 4$ et donc $(q\ell - q + c) = 10$. Comme $G^2(M) = 0.18$ et $G^2(M^*) = .16$, on obtient $\text{AIC}(M) = 18.18$ et $\text{AIC}(M^*) = 20.16$. De même, on trouve $\text{BIC}(M) = 41.63$ et $\text{BIC}(M^*) = 46.21$. Les deux critères indiquent que l'arbre M plus simple est préférable à l'arbre M^* . Le gain en qualité d'ajustement de M^* n'est pas assez important pour justifier l'accroissement de la complexité. Remarquons que du point de vue de ces critères d'information, l'arbre induit est supérieur tant au modèle d'indépendance ($\text{AIC}(I) = 32.55$, $\text{BIC}(I) = 50.78$) qu'au modèle saturé ($\text{AIC}(S) = 24$, $\text{BIC}(S) = 55.26$).

5. Illustration

Nous illustrons ici les enseignements apportés par les critères de qualité d'ajustement proposés sur un exemple concret. On considère pour cela les données relatives aux 762 étudiants qui ont commencé leur première année d'études à la Faculté des sciences économiques et sociales de Genève en 1998. Il s'agit de données administratives réunies par [PET 01]. On rapporte quelques résultats d'une analyse visant à

évaluer les chances de respectivement réussir, redoubler ou être éliminé à la fin de la première année d'études selon les caractéristiques personnelles portant notamment

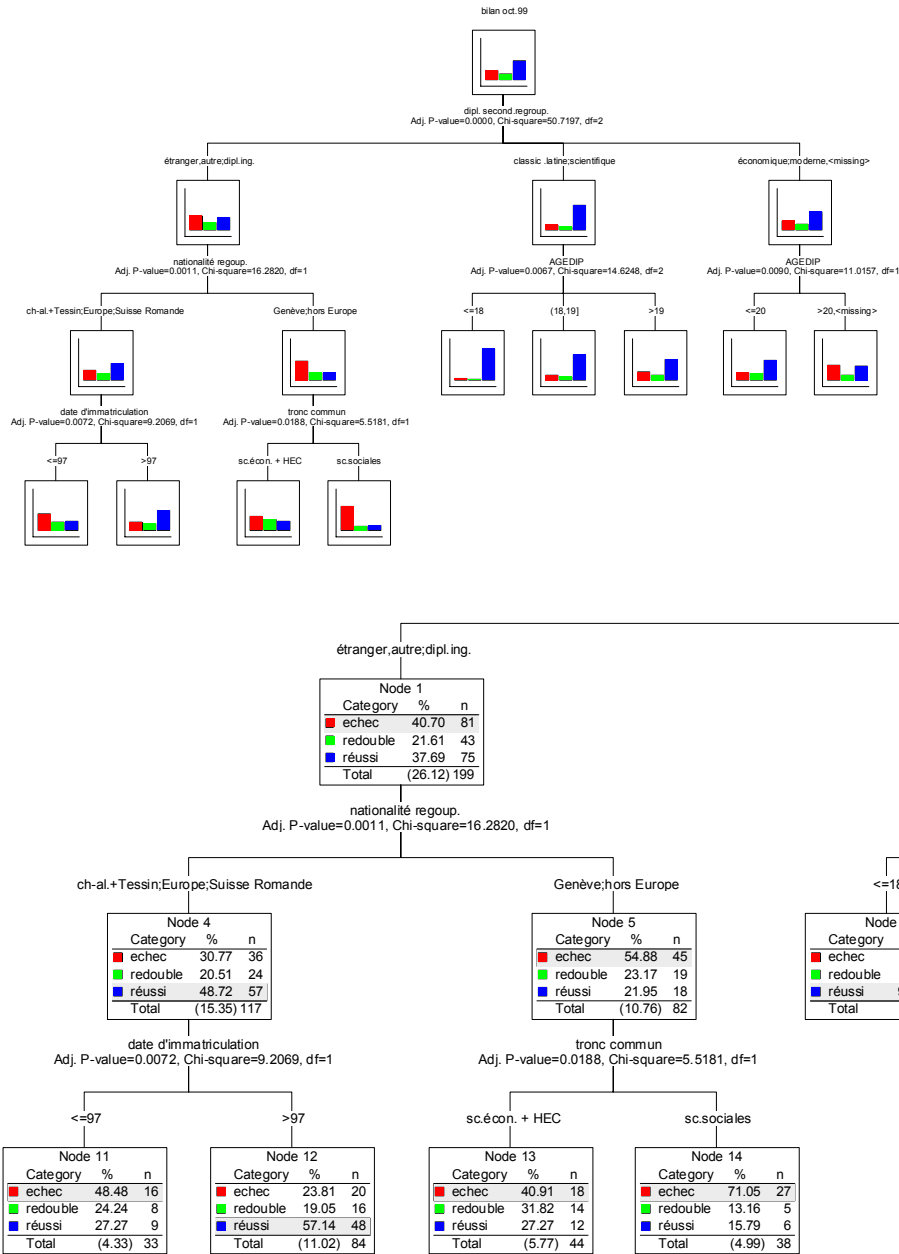


Figure 4. Bilan après une année en SES : arbre CHAID et détail de la branche gauche

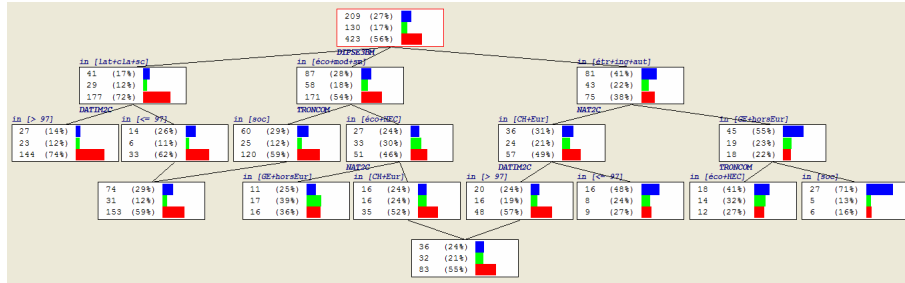


Figure 5. Graphe induit avec Sipina

sur l'origine et le cursus scolaire. La figure 4 montre l'arbre obtenu avec la procédure CHAID [KAS 80] implémentée dans Answer Tree [SPS 01]. Parmi une trentaine de prédicteurs potentiels, CHAID en a sélectionné 5 dont deux quantitatifs, l'année d'immatriculation à l'université et l'âge à l'obtention du diplôme de l'école secondaire. Les cinq variables avec les discrétisations et regroupements de modalités proposés par CHAID sont le type de diplôme secondaire (3 modalités), l'âge de son obtention (4), la date d'immatriculation (2), le tronc commun choisi (2) et la nationalité (2). La table cible **T** définie par ces variables contient 88 colonnes. Elle a 3 lignes correspondant aux 3 situations possibles de l'étudiant après sa première année d'étude.

Le tableau 3 donne la taille de la partition, la déviance G^2 avec ses degrés de liberté et son degré de signification et les critères AIC et BIC. Ces valeurs peuvent être comparées à celles de plusieurs variantes. CHAID2 est CHAID sans l'éclatement du sommet 4 ($nationa \notin \{GE, \text{hors Europe}\}$) et CHAID3 sans l'éclatement des sommets 4 et 5 ($nationa \in \{GE, \text{hors Europe}\}$). Le modèle Sipina correspond au graphe de la figure 5 obtenu avec la procédure Sipina [Sip 00][ZIG 00] qui, comme on peut le voir, autorise également des fusions de sommets. On donne également les valeurs trouvées pour les partitions qui donnent respectivement les plus petits AIC et BIC. Enfin, le

Tableau 3. SES 98 : qualités d'ajustement d'un choix de modèles

Modèle	q	d	G^2	$\text{sig}(G^2)$	AIC	BIC
Saturé	88	0	0	1	528	1751.9
Meilleur AIC	14	148	17.4	1	249.4	787.2
CHAID	9	158	177.9	0.133	390.0	881.3
CHAID2	8	160	187.4	0.068	395.4	877.5
CHAID3	7	162	195.2	0.038	399.2	872.1
Sipina	7	162	185.8	0.097	389.8	862.6
Meilleur BIC	6	164	75.2	1	275.2	738.8
Indépendance	1	174	295.1	0.000	475.8	892.3

CHAID2 : CHAID sans éclatement *datimma* du sommet 4 ($nationa \neq GE, \text{hors Europe}$)

CHAID3 : CHAID2 sans éclatement *troncom* du sommet 5 ($nationa = GE, \text{hors Europe}$)

Tableau 4. SES 98 : mesures de type R^2 pour un choix de modèles

Modèle	proportion de réduction d'entropie			relativement au modèle saturé			pseudo R^2_{ajust}
	$\hat{\tau}_{M I}$	$\hat{u}_{M I}$	$\hat{\lambda}_{M I}$	$\hat{\tau}_{M I}$	$\hat{u}_{M I}$	$\hat{\lambda}_{M I}$	
Saturé	0.193	0.197	0.183	1	1	1	1
Meilleur AIC	0.159	0.185	0.179	0.824	0.939	0.978	.941
CHAID	0.094	0.078	0.109	0.487	0.396	0.596	.336
CHAID2	0.086	0.072	0.088	0.446	0.365	0.481	.309
CHAID3	0.08	0.067	0.088	0.415	0.340	0.481	.289
Sipina	0.087	0.073	0.103	0.451	0.371	0.563	.324
Meilleur BIC	0.149	0.147	0.142	0.772	0.746	0.776	.745
Indépendance	0	0	0	0	0	0	0

modèle saturé correspond à la partition la plus fine et le modèle d'indépendance au cas où tous les profils sont regroupés en seul groupe.

On constate tout d'abord qu'à l'exception du modèle d'indépendance et de CHAID3 avec un degré de signification légèrement inférieur à 5%, tous les modèles reproduisent de façon satisfaisante la table T. La simplification de l'arbre CHAID en CHAID2 ou CHAID3 se traduit comme attendu par une détérioration du G^2 . Les écarts sont $G^2(\text{CHAID2}|\text{CHAID}) = 9.5$ et $G^2(\text{CHAID3}|\text{CHAID}) = 17.3$ qui pour un gain de respectivement 2 et 4 degrés de liberté sont clairement significatifs, ce qui valide statistiquement l'éclatement des nœuds 4 et 5. Les différences de G^2 avec les autres modèles qui ne sont pas des sous graphes de l'arbre CHAID ne peuvent être testés. On peut par contre comparer avec les AIC ou BIC de ces modèles. On remarque tout d'abord que CHAID et ses deux variantes ont des AIC et BIC très voisins, sensiblement meilleurs que ceux du modèle saturé et dans une moindre mesure que ceux du modèle d'indépendance. La partition générée par Sipina obtient un AIC équivalent au modèle CHAID, mais son BIC est inférieur à celui des 3 modèles CHAID. L'écart supérieur à 10 traduit, selon l'échelle postulée par Raftery [RAF 95], une supériorité très forte de cette partition. On peut noter toutefois, que les valeurs des AIC et BIC obtenues pour le graphe Sipina restent très nettement supérieures aux valeurs optimales possibles avec les attributs retenus.

Le tableau 4 récapitule les mesures de type R^2 . Le $\hat{\lambda}_{M|I}$ qui donne la proportion de réduction du taux d'erreur sur données d'apprentissage est donné pour comparaison avec les proportions de réduction d'entropie. Les proportions maximales de réduction d'entropie par rapport au modèle d'indépendance sont évidemment obtenues avec la partition la plus fine, c'est-à-dire le modèle saturé. On note que ces maxima sont inférieurs à 1. La part de cette réduction maximale réalisée par chaque modèle est également donnée. On voit que ces dernières valeurs sont très similaires au pseudo R^2 ajusté. Elles nous indiquent par exemple, que les modèles CHAID et Sipina, malgré un ajustement satisfaisant, ne captent qu'environ 1/3 du potentiel de réduction d'entropie possible avec les prédicteurs retenus. Les meilleures partitions du point de vue tant du BIC que de l'AIC font nettement mieux de ce point de vue. On notera cependant que

les règles (non données ici) qui caractérisent les classes des partitions optimales sont très complexes et donc difficiles à décrire.

6. Conclusion

Cet article aborde la question de la qualité de l'ajustement des arbres d'induction. Il s'agit d'un aspect peu discuté dans la littérature sur l'extraction de connaissances alors même que la qualité d'ajustement fait partie des outils classiques d'évaluation de modèles en statistique. La qualité d'ajustement fournit des indications complémentaires aux indicateurs de qualité traditionnellement utilisés pour les arbres d'induction en permettant, en particulier, d'évaluer la pertinence statistique d'un arbre induit.

Concrètement, nous avons montré, en introduisant les notions d'arbre saturé et d'arbre étendu, comment adapter aux arbres d'induction les statistiques du khi-2 de Pearson et du rapport de vraisemblance utilisés dans le cadre de la modélisation de tables de contingence. Nous avons également considéré la question de la comparaison de modèles pour laquelle la différence des statistiques G^2 du rapport de vraisemblance permet de tester la significativité statistique du gain d'information d'un modèle par rapport à un modèle de référence. Enfin, nous avons vu que l'on pouvait exploiter également les critères d'information AIC et BIC pour guider le choix entre modèles de complexité variable.

Ce travail avait pour objectif de montrer comment appliquer des critères statistiques bien établis aux arbres d'induction. Il reste évidemment encore beaucoup à faire. D'une part, il convient de mettre en œuvre dans des cas concrets les statistiques discutées et en particulier de les implémenter dans une procédure de construction d'arbres d'induction.

L'approche retenue dans cet article qui s'appuie notamment sur les concepts d'arbre saturé et d'extension maximale de l'arbre, s'applique lorsque le croisement de toutes les modalités des attributs donne lieu à un nombre raisonnable de catégories. Dans le cas particulier de variables quantitatives continues, il y a lieu de discrétiser les valeurs. La difficulté tient ici au fait que, dans les arbres d'induction, la discrétisation ne se fait en règle générale pas a priori mais est déterminée en cours de processus de façon à optimiser la discrimination entre classes.

Une approche possible est de retenir la discrétisation des variables continues définie par l'ensemble des seuils utilisés par le graphe induit. Les variables continues étant ainsi rendues catégorielles, la construction de l'arbre saturé et de l'arbre étendu devient possible et la démarche précédente s'applique. Les résultats sont alors conditionnels à la discrétisation retenue, ce qui en limite évidemment la portée. On peut songer à d'autres approches prenant en particulier en compte le fait que les seuils de discrétisation sont également des paramètres du modèle. Ceci mérite cependant une réflexion approfondie qui dépasse le cadre de cet article

Enfin, de façon plus générale, il reste encore beaucoup de questions ouvertes sur la pertinence statistique des arbres d'induction. Par exemple, la mesure de la fiabilité des estimations des paramètres du modèle de reconstruction (1) issu de l'arbre et celle de la stabilité de l'arbre induit sont à nos yeux essentielles pour apprécier la confiance à accorder à un arbre.

7. Bibliographie

- [AGR 90] AGRESTI A., *Categorical Data Analysis*, Wiley, New York, 1990.
- [AKA 73] AKAIKE H., « Information Theory and an Extension of the Maximum Likelihood Principle », in PETROX B. N., CASIKI F., Eds., *Second International Symposium on Information Theory*, page 267, Akademiai Kiado, Budapest, 1973.
- [BIS 75] BISHOP Y. M. M., FIENBERG S. E., HOLLAND P. W., *Discrete Multivariate Analysis*, MIT Press, Cambridge MA, 1975.
- [BRE 84] BREIMAN L., FRIEDMAN J. H., OLSHEN R. A., STONE C. J., *Classification And Regression Trees*, Wadsworth International Group, Belmont, CA, 1984.
- [GEU 02] GEURTS P., « Contributions to Decision Tree Induction : Bias/Variance Tradeoff and Time Series Classification », PhD Thesis, Université de Liège, Faculté des Sciences Appliquées, Liège, May 2002.
- [GOO 54] GOODMAN L. A., KRUSKAL W. H., « Measures of association for cross classifications », *Journal of the American Statistical Association*, vol. 49, 1954, p. 732–764.
- [GOO 72] GOODMAN L. A., KRUSKAL W. H., « Measures of association for cross classifications IV : simplification of asymptotic variances », *Journal of the American Statistical Association*, vol. 67, 1972, p. 415–421.
- [KAS 80] KASS G. V., « An exploratory technique for investigating large quantities of categorical data », *Applied Statistics*, vol. 29, n° 2, 1980, p. 119–127.
- [KAS 95] KASS R. E., RAFTERY A. E., « Bayes Factors », *Journal of the American Statistical Association*, vol. 90, n° 430, 1995, p. 773–795.
- [LIG 71] LIGHT R. J., MARGOLIN B. H., « An Analysis of Variance for Categorical Data », *Journal of the American Statistical Association*, vol. 66, n° 335, 1971, p. 534–544.
- [OLS 95] OLSZAK M., RITSCHARD G., « The behaviour of nominal and ordinal partial association measures », *The Statistician*, vol. 44, n° 2, 1995, p. 195–212.
- [PET 01] PETROFF C., BETTEX A.-M., KORFFY A., « Itinéraires d'étudiants à la Faculté des Sciences économiques et sociales : le premier cycle », rapport, Juin 2001, Université de Genève, Faculté SES.
- [POW 00] POWERS D. A., XIE Y., *Statistical Methods for Categorical Data Analysis*, Academic Press, San Diego, CA, 2000.
- [QUI 93] QUINLAN J. R., *C4.5 : Programs for Machine Learning*, Morgan Kaufmann, San Mateo, 1993.
- [RAF 95] RAFTERY A. E., « Bayesian Model Selection in Social Research », in MARSDEN P., Ed., *Sociological Methodology*, p. 111-163, The American Sociological Association, Washington, DC, 1995.

- [RAK 98] RAKOTOMALALA R., ZIGHED D. A., « Mesures PRE dans les graphes d'induction : une approche statistique de l'arbitrage généralité-précision », in RITSCHARD G., BERCHTOLD A., DUC F., ZIGHED D. A., Eds., *Apprentissage : des principes naturels aux méthodes artificielles*, p. 37–60, Hermes Science Publications, Paris, 1998.
- [RIT 03] RITSCHARD G., ZIGHED D. A., « Modélisation de tables de contingences par arbres d'induction », *RTSI Extraction des connaissances et apprentissage*, vol. 17, n° 1–3, 2003, p. 381–392.
- [SCH 78] SCHWARZ G., « Estimating the Dimension of a Model », *The Annals of Statistics*, vol. 6, 1978, p. 461–464.
- [Sip 00] SIPINA FOR WINDOWS V2.5, <http://eric.univ-lyon2.fr>, 2000, Logiciel.
- [SPS 01] SPSS, Ed., *Answer Tree 3.0 User's Guide*, SPSS Inc., Chicago, 2001.
- [THE 67] THEIL H., *Economics and Information Theory*, North-Holland, Amsterdam, 1967.
- [THE 70] THEIL H., « On the Estimation of Relationships Involving Qualitative Variables », *American Journal of Sociology*, vol. 76, 1970, p. 103–154.
- [ZIG 00] ZIGHED D. A., RAKOTOMALALA R., *Graphes d'induction : apprentissage et data mining*, Hermes Science Publications, Paris, 2000.