

# Critères d'information pour la sélection de variables

Mohamed Abadi<sup>1</sup>, Olivier Alata<sup>1</sup>, Christian Olivier<sup>1</sup>, Enguerran Grandchamp<sup>2</sup> et  
Majdi Khoudeir<sup>1</sup>

<sup>1</sup> Université de Poitiers,  
Institut XLIM-SIC, UMR CNRS 6172,  
BP. 30179, 86962 Futuroscope-Chasseneuil Cedex, France.  
{abadi,alata,olivier,khoudeir}@sic.univ-poitiers.fr

<sup>2</sup> Université des Antilles et de la Guyane,  
Laboratoire LAMIA,  
Campus de Fouillole, 97157, Pointe-à-Pitre, Guadeloupe, France  
egradch@univ-ag.fr

**Résumé** Cet article étudie les potentialités des critères d'information pour la sélection des variables. Ces critères sont intégrés dans des schémas de parcours de sous-ensembles candidats afin de sélectionner ceux qui les minimisent. Le choix d'un sous-ensemble dépend de la qualité de l'approximation des densités de probabilité de ces variables. A cet effet, nous avons adopté une approche exploitant les histogrammes optimisés à l'aide des principes en lien avec le codage arithmétique adaptatif. Des tests sur des données simulées et de référence sont réalisés. Ils montrent par le biais des taux de bonne classification, déterminés en appliquant plusieurs algorithmes de classification, l'importance d'un tel outil et sa capacité à choisir les sous-ensembles qui caractérisent mieux les classes vis-à-vis de données.

**Mots clés** Critères d'information, sélection de variables, estimation et sélection des histogrammes, schémas de parcours de variables.

**Abstract** This paper introduces the information criteria in feature selection framework. Information criteria are integrated in feature selection scheme to select subset candidates. The accuracy of the proposed approach is based on the quality of the probability density approximations of these features. They are obtained using histograms optimized thanks to the adaptive arithmetic coding principles. Tests on simulated data and references are made. Multiple classifiers are used. The correct classification rate shows, the importance of this tools and its ability to select a best subsets. Generally, the subsets produce a good characterization of classes which the data belong.

**Key words** Information criteria, feature selection, histogram estimation and selection, feature scheme search.

## 1 Introduction

La sélection de variables connue aussi sous le nom de sélection d'attributs est considérée comme l'une des étapes clés dans plusieurs domaines (ex. l'apprentissage automatique, la modélisation de données,

le data-mining, etc. [1, 2, 3]). Elle permet de réduire la dimensionnalité des données en identifiant un ou plusieurs sous-ensembles appartenant à l'ensemble d'origine qui est composé par toutes les variables. Par exemple soient  $X = \{x_i, i = 1, \dots, M\}$  une observation à  $M$  variables qui représente la classe  $k$ . L'objectif est donc de trouver parmi les  $M$  variables les  $m$  variables ( $m < M$ ) qui caractérisent au mieux  $k$ . Dans ce cas, le nombre total de sous-ensembles qui peuvent être construit à partir de  $M$  variables est  $2^M - 1$  et ceux dont la dimension est inférieure ou égale à  $m$  est  $\sum_{i=1}^m \binom{M}{i}$ . Plusieurs algorithmes ont été développés pour construire les sous-ensembles car une recherche exhaustive sur les  $M$  variables est très coûteuse en temps de calcul. Pour plus de détails, le lecteur intéressé peut se référer à [4]. Nous nous intéressons dans cette étude aux algorithmes basés sur des approches séquentielles et plus particulièrement aux schémas de sélection progressive (SFS pour Sequential Forward Selection) et rétrograde (SBS pour Séquential Backward Selection). L'avantage est qu'ils sont faciles à implémenter, indépendants des algorithmes de classification et des critères de décision. Donc ces schémas offrent la possibilité de comparer aisément plusieurs critères entre eux. L'important est de trouver les sous-ensembles qui traduisent au mieux l'information recherchée. Ceci met en avant l'importance du choix du critère de sélection.

Différents critères ont été développés pour évaluer les sous-ensembles candidats. Ils explorent des mesures de distance [5], des mesures statistiques [6] ou plus récemment des mesures probabilistes basées sur l'estimation de l'information mutuelle [7]. Dans cet article, nous proposons d'étudier le potentiel des critères d'information basés sur un calcul du maximum de vraisemblance (MV) pénalisé pour la sélection de variables [8]. Ces critères sont intégrés dans les schémas SFS et SBS. La détermination du MV exige l'estimation de la distribution des variables. En effet, nous adoptons ici le codage arithmétique pour estimer la densité de probabilité et la partition optimale des variables [9].

## 2 Critères d'information

Dans ce travail nous proposons l'adaptation et l'intégration des critères d'information AIC [8], BIC [10] et  $\varphi_\beta$  [11] (resp. pour Akaike, Bayes Information Criterion et phi-beta) dans un schéma de sélection de variables. Rappelons que la forme générale d'un critère d'information IC pour un modèle paramétré  $Y$  s'écrit sous la forme suivante :

$$IC(Y) = -\log(MV) + |Y|\alpha(n) \quad (1)$$

où  $|Y|$  est le nombre de paramètres libres qui croît en fonction de la complexité du sous-ensemble et  $\alpha(n)$  est une fonction pénalisante qui varie en fonction du critère utilisé (respectivement les  $\alpha(n)$  de AIC, BIC et  $\varphi_\beta$  sont égaux à 2,  $\log n$  et  $n^\beta \log \log n$  avec  $0 < \beta < 1$  pour assurer de bonnes propriétés asymptotiques). La quantité  $|Y|\alpha(n)$  permet de régulariser le comportement du MV car il est connu que ce dernier sur-paramétrise. Ainsi la minimisation du critère permet d'avoir un compromis entre l'adéquation aux données et la complexité des sous-ensembles candidats. En vue de résoudre un problème de sélection de variables à l'aide des ICs, nous nous plaçons dans le cas suivant : soient  $X$  un ensemble à  $M$  variables et  $Y \subseteq X$ , sous-ensemble de  $X$  à  $m$  variables. On dispose de  $n$  observations. Chaque observation est assignée à une classe  $k$ ,  $k = 1, \dots, K$ .  $K$  est le nombre total de classes. La probabilité jointe d'avoir la classe  $k$  et l'ensemble  $X$  peut être formulée par  $P(k, X) = \frac{P(k, Y) \cdot P(X)}{P(Y)}$  dans le cas où la classification devient indépendante du sous-ensemble  $X \setminus Y$  conditionnellement à  $Y$  [8]. Nous supposons que  $P(Y) = 1$  pour  $Y = \emptyset$  (sous-ensemble vide). D'après [8], Les ICs sont ainsi définis par

$$IC(Y) = -\sum_{k=1}^K \sum_{c \in C_Y} n_{k,c} \log \left\{ \frac{n_{k,c}}{n(Y)} \right\} + (|C_{k,Y}| - |C_Y|)\alpha(n) \quad (2)$$

78 où  $|C_{k,Y}|$  et  $|C_Y|$  dénotent respectivement le nombre de zones (appelées ensuite « boîtes ») de  
 79 l'histogramme multidimensionnel d'occurrence non nul en lien avec le sous-ensemble  $Y$ , d'une part en  
 80 prenant compte des différentes classes et d'autre part sans tenir compte de la classification. Nous  
 81 supposons que  $0 \times \log 0 = 0$  et  $n(Y) = n$ .  $n_{k,c}$  est donc le nombre d'observations se trouvant dans la  
 82 boîte  $c$  de l'histogramme multidimensionnel. L'obtention de l'histogramme multidimensionnel est  
 83 réalisée en recherchant au préalable un histogramme optimal pour chaque variable.

### 84 3 Estimation de la distribution des variables

85 L'estimation optimale de la densité de probabilité (ddp) est un problème difficile. Cette estimation  
 86 s'appuie sur un ensemble fini d'observations  $v^n = (v_1, \dots, v_n)$ , réalisation de  $n$  variables aléatoires  
 87 indépendantes et identiquement distribuées (iid)  $(V_1, \dots, V_n)$  de ddp  $f$  qui est donc à estimer à partir de  
 88  $v^n$ . Nous citons ici les méthodes non paramétriques les plus populaires : les méthodes à noyau et les  
 89 méthodes par histogramme. Nous nous intéressons particulièrement à l'estimation de la ddp par  
 90 histogramme en utilisant le codage arithmétique adaptatif [9].

91 Soient  $f$  une ddp inconnue définie sur l'intervalle  $I = [a, b] \subset \mathbb{R}$  et  $v^n$  une observation. Soit  $P =$   
 92  $(I_j)_{j \in 1, \dots, m}$  une partition de  $I$  à  $m \in \mathbb{N}^*$  intervalles. Une estimation de  $f$  au sens du MV à partir de  $v^n$   
 93 peut être obtenue en maximisant l'expression suivante :

$$94 \quad \hat{f}_P = \sum_{j=1}^m \frac{n_j}{n|I_j|} \mathbb{1}_{I_j} \quad (3)$$

95 où  $n_j$  est le nombre de valeurs issues de  $v^n$  tombant dans l'intervalle  $I_j$ ,  $|I_j|$  est la longueur de  $I_j$  et  $\mathbb{1}_{I_j}$   
 96 dénote la fonction indicatrice de l'intervalle  $I_j$ . La difficulté d'une telle procédure d'estimation est de  
 97 trouver la partition  $P$  qui estime au mieux  $f$  par  $\hat{f}_P$ . Pour cela, Coq et al. [9] ont élaboré une technique  
 98 basé sur les critères d'information et la longueur du codage sans perte des données  $v^n$ . Cette  
 99 technique est basée sur deux types de codage : le codage arithmétique et le codage de précision. Nous  
 100 rappelons ici le critère résultant de ces développements.

$$101 \quad CRIT(v^n | P) = - \sum_{j=1}^m n_j \log \frac{n_j}{n|I_j|} + \frac{(m-1)}{2} \alpha(n) - n \log r \quad (4)$$

102 Le principe MDL (pour Minimum Description Length) suggère de choisir la partition qui minimise ce  
 103 critère pour approcher  $f$  par  $\hat{f}_P$  de la manière suivante :

$$104 \quad \hat{P} = \operatorname{argmin}\{CRIT(v^n, P), P \in P_{max}\} \quad (5)$$

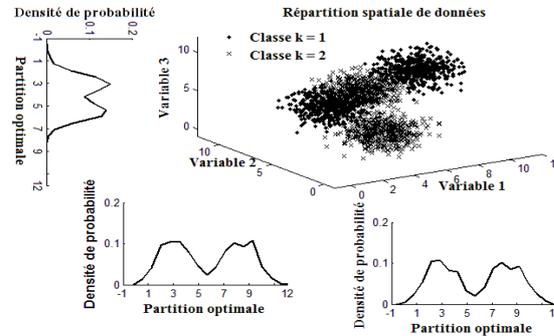
105 où  $P_{max}$  est la partition maximale de  $I$  à  $R$  intervalles tous de longueur  $r$  définie à priori.

106 L'obtention d'une telle partition et d'une telle densité pour chaque variable à l'aide de la  
 107 programmation dynamique afin de ne pas tester toutes les partitions possibles constitue le point de  
 108 départ pour calculer l'équation (2).

### 109 4 Expérimentations

110 Notre approche de sélection de variables est basée sur les critères d'information définis par l'équation  
 111 (2). Elle est évaluée et comparée d'abord sur des données simulées qui suivent une loi normale  
 112  $\mathcal{N}(\mu, \sigma)$  où le nombre de classes  $K$  est égal à 2, le nombre de variables  $M = 3$  et  $n = 500$  par classe

113 (figure 1), ensuite sur le jeu de données WINE<sup>1</sup> dont les caractéristiques sont  $K = 3$ ,  $n = 178$ , et  
 114  $M = 13$ ). Pour ce jeu de données, nous avons utilisé la validation croisée (CV pour cross validation)  
 115 pour construire les ensembles d'apprentissage et de test ( $10 - CV$ ). Cette opération est répétée  $N$ -fois  
 116 pour obtenir des résultats statistiquement corrects ( $N = 100$ ).



117

118 **Figure 1** : Répartitions spatiales des données et ddp. Les variables suivent une loi normale  $\mathcal{N}(\mu, \sigma = 1)$  sur  
 119 chaque axe. Le tableau 1 détaille les valeurs de la moyenne  $\mu$  pour chaque variable et chaque classe ( $K = 2$ ).

120 **Tableau 1** : Valeurs de moyenne  $\mu$  présent pour chaque variable en fonction de chaque classe.

Numéro classe	Variable 1	Variable 2	Variable 3
$k = 1$	(2,10)	(4,8)	(6,6)
$k = 2$	(4,8)	(2,10)	(3,3)

121 A partir des ddp estimées en utilisant l'équation (4), nous avons calculé les ICs de l'équation (2) en  
 122 utilisant les schémas de parcours SFS et SBS. Le but est de quantifier la qualité des variables afin de  
 123 sélectionner celles qui décrivent au mieux les données. Le tableau 2 montre les sous-ensembles  
 124 sélectionnés à chaque itération par le critère d'information BIC dans le cas de la figure 1. Ces résultats  
 125 sont comparés au critère de Wilk's tel qu'il est défini dans [12]. En analysant uniquement les ddp, on  
 126 observe que le sous-ensemble formé uniquement par la variable 3 décrit bien les données. Un sous-  
 127 ensemble contenant cette variable ne peut que mieux représenter ces données. Cependant, le sous-  
 128 ensemble  $\{1,2\}$  sépare presque parfaitement les données. Le tableau 2 montre aussi la différence du  
 129 sous-ensemble obtenu à l'itération 2 par le schéma SBS entre les critères BIC ( $Y_{SBS} = \{1,2\}$ ) et Wilk's  
 130 ( $Y_{SBS} = \{2,3\}$ ). Cependant, ces critères, pour le schéma SFS, renvoi le même sous ensemble ( $Y_{SFS} =$   
 131  $\{2,3\}$ ). D'où l'intérêt du schéma SBS en utilisant le critère BIC.

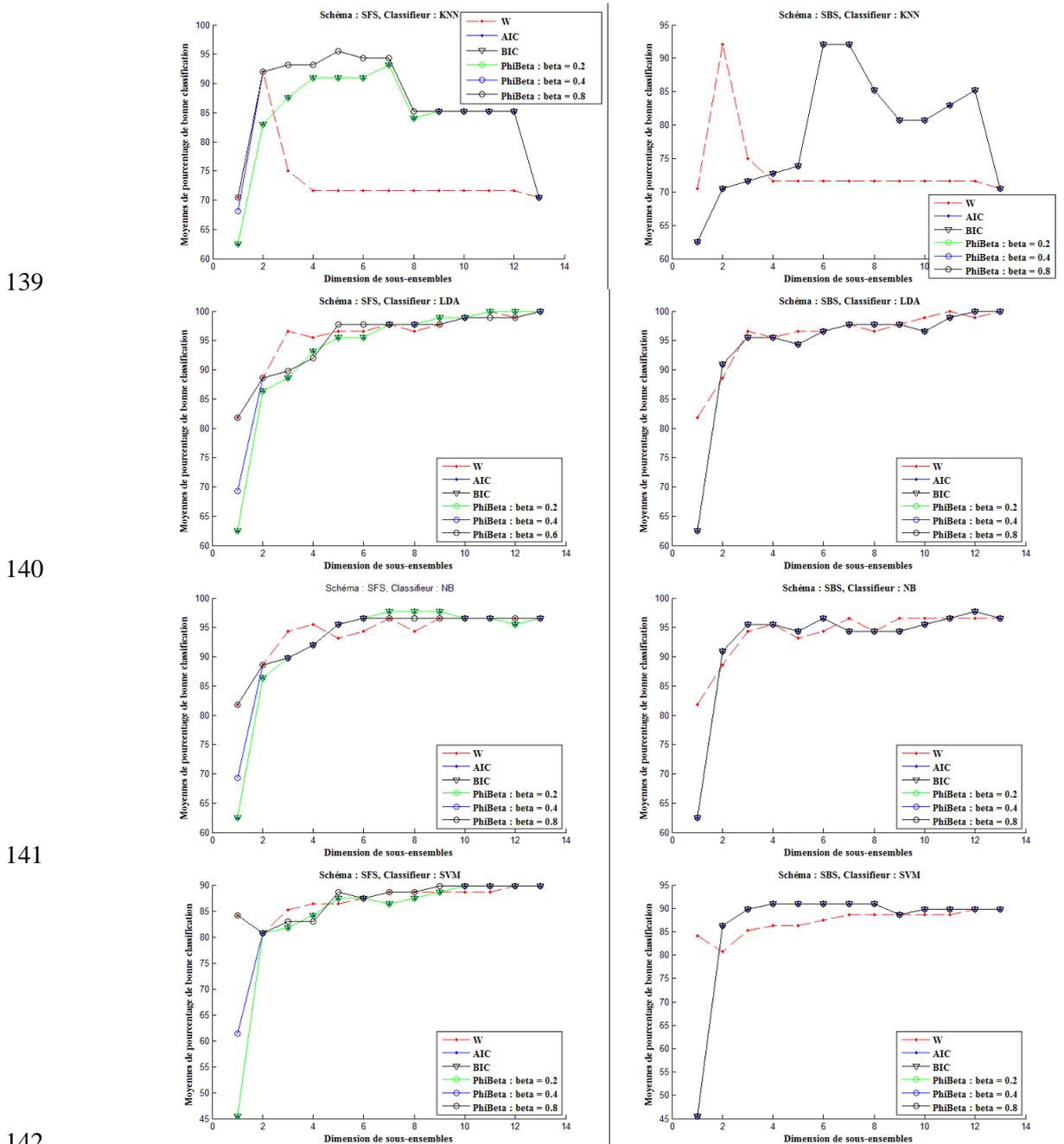
132

**Tableau 2** : Résultats de sous-ensembles sélectionnés.

	SFS		SBS	
	BIC	Wilk's	BIC	Wilk's
<i>Itération 1</i>	{3}	{3}	{1,2,3}	{1,2,3}
<i>Itération 2</i>	{2,3}	{2,3}	{1,2}	{2,3}
<i>Itération 3</i>	{1,2,3}	{1,2,3}	{2}	{3}

<sup>1</sup> Base de données UCI Repository : <http://archive.ics.uci.edu/ml/datasets.html>

133 La figure 2 ci-dessous présente une comparaison quantitative entre les différents sous-ensembles  
 134 obtenus par les trois critères d'information et le critère de Wilk's. Ces résultats montrent les  
 135 pourcentages moyens de taux de bonne classification obtenus par les algorithmes de classification.  
 136 Quatre algorithmes de classification ont été appliqués pour mesurer la qualité des résultats et les  
 137 performances des schémas SFS et SBS à savoir  $k$ -NN ( $k = 1$ ,  $k$ -Nearest Neighbor), NB (Naive  
 138 Bayes), LDA (Linear Discriminant Analysis) et SVM (Support Vector Machine).



143 **Figure 2 :** Comparaison de taux de bonne classification entre les critères d'information et celui de Wilk's pour  
 144 le jeu de données WINE. Par lignes classifieurs  $k$ -NN, LDA, NB et SVM. Par colonnes schémas SFS et SBS.

145 La figure 2 montre la supériorité des critères d'information par rapport au critère de Wilk's. Les  
 146 courbes issues des critères d'information se confondent dans la plupart de cas. Cela signifie que ces  
 147 critères ont sélectionné les mêmes sous-ensembles pour les différentes itérations. Pour le critère  $\varphi_\beta$   
 148 nous avons choisi  $\beta = 0.2, 0.4$  et  $0.8$ . Pour ces valeurs de  $\beta$ , nous constatons que le taux de bonne  
 149 classification est le plus élevés. L'utilisation de différents algorithmes de classification affecte peu  
 150 notre approche. Elle converge mieux vers la solution optimale obtenue par une recherche exhaustive.

## 151 5 Conclusion

152 Nous avons présenté une approche de sélection de variables basée sur les critères d'information. Des  
 153 estimateurs continus par morceaux utilisant le codage arithmétique adaptatif sont employés pour  
 154 approcher la ddp et définir ainsi la partition optimale pour chaque variable. Les tests réalisés sur des  
 155 données simulées et références et les comparaisons avec le critère de Wilk's ont montré la bonne  
 156 tenue de notre approche pour différents algorithmes de classification.

## 157 Remerciements

158 Les auteurs remercient les institutions européennes à travers le programme INTERREG IV (projet  
 159 CESAR volet II) dont l'Université des Antilles et de la Guyane est partenaire et la région Poitou-  
 160 Charentes d'avoir financé ce projet.

## 161 Références

- 162 [1] A.K. Jain, R.P.W. Duin, and J. Mao. Statistical pattern recognition: A review. *IEEETrans. PAMI*, 22(1)  
 163 :4-37, 2000.
- 164 [2] S. Mitraetal. Datamining in soft computing frame work: A survey. *IEEETrans.NN*, 13(1) :3-14,2002.
- 165 [3] H. Liuand and L.Yu. Toward integrating FS algorithms for classification and clustering. *IEEETrans. on*  
 166 *KDE*, 17(4) :491-502, 2005.
- 167 [4] Y. Sun, S. Todorovic, and S. Goodison. Local-Learning-Based Feature Selection for High-Dimensional  
 168 Data Analysis. *TPAMI*, 32,(9) :1610-1626, 2010.
- 169 [5] J. Liang, Su Yang and A-C. Winstanley. Invariant optimal feature selection: A distance discriminant and  
 170 feature ranking based solution. *Pattern Recognition*, 41(5) :1429-1439, 2008.
- 171 [6] G. Qu, S. Hariri and M. Yousif. A new dependency and correlation analysis for features. *IEEE*  
 172 *Transactions on Knowledge and Data Engineering*, 17(9) :1199-1207, 2005.
- 173 [7] H. Liu, J. Sun, L. Liuand and H. Zhang. Feature selection with dynamic mutual information. *Pattern*  
 174 *Recognition*, 42(7) :1330-1339, 2009.
- 175 [8] Y. Sakamoto and H. Akaike. Analysis of cross classified data by AIC. *Ann. Inst. Statist. Math*, 30(B)  
 176 :185-197, 1978.
- 177 [9] G. Coq, O. Alata, Y. Pousset, X. Li and C. Olivier. Law recognition via histogram-based estimation. *IEEE*  
 178 *ICASSP, Taïpei (Taïwan)*, 3425-3428, avril 2009.
- 179 [10] T. Mary-Huard and E. Lebarbier. Une introduction au critère BIC : fondements théoriques et applications.  
 180 *Journal de la société française de Statistiques*, 147(1) :39-57, 2006.
- 181 [11] D.A. El-Matouat and M. Hallin. Order selection, stochastic complexity and Kullback-Leibler information.  
 182 115 :291-299, 1996.
- 183 [12] A. Porebski, N. Vandenbroucke and L. Macaire. Comparison of feature selection schemes for color texture  
 184 classification. *IPTA, Paris (France)*, july 2010.