

# Arbres de décision k-means

Célia Châtel, François Brucker, Pascal Pr ea

Aix-Marseille Universit e - LIF -  cole Centrale Marseille  
{celia.chatel, francois.brucker, pascal.prea}@lif.univ-mrs.fr

**Mots-cl es :** *classification, arbre de d cision, k-means, syst me totalement  quilibr *

## 1 Introduction

Les arbres de d cision (Breiman [2]) sont un mod le tr s utilis  en apprentissage automatique supervis . Dans ce travail, nous appliquons des id es de classification non supervis e (clustering)   la construction d'arbres semblables   des arbres de d cision. Les arbres sont ainsi construits en utilisant uniquement des crit res de distance et non de puret  des classes produites. Nous pr sentons les arbres de d cision obtenus et un crit re de performance.

## 2 M thode

L'Algorithme (1) d taille la construction d'un arbre de k-means (KMDT). Le syst me de classes obtenu est une hi rarchie et reprend le principe des k-means hi rarchiques (B cker [1]). L'utilisation du mod le ainsi construit pour une t che de pr diction consiste   faire descendre les exemples dans la structure   la mani re d'un arbre de d cision : l'exemple descend dans le noeud fils dont il est le plus proche du centre de gravit .

On calcule au maximum  $n$  k-means, mais en pratique beaucoup moins. Un arbre de d cision classique contient moins de noeuds mais doit   chaque noeud optimiser le d coupage selon toutes les variables et peut tester de nombreux d coupages en cas de variables r elles.

```
Cr er un noeud racine contenant tous les exemples
Cr er une file  $\mathcal{F}$  contenant la racine
tant que la file  $\mathcal{F}$  est non vide faire
  | D piler de  $\mathcal{F}$  un noeud
  | si tous les exemples que le noeud contient ne sont pas de la m me classe alors
  | | Calculer un k-means avec  $k = 2$ 
  | | Cr er un noeud fils du noeud courant pour chaque cluster obtenu
  | | Empiler les noeuds fils dans  $\mathcal{F}$ 
  | fin
fin
```

**Algorithme 1 :** Algorithme de construction d'un arbre de k-means

Les mod les les plus efficaces bas s sur des arbres de d cision classiques sont les for ts d'arbres al atoires (Breiman [3]) construits   partir de l'id e de sous-espaces al atoires (Ho [5]). Les KMDT peuvent  tre rendus al atoires gr ce   la m me id e en appliquant la division de chaque noeud par le k-means   partir d'un nombre limit  de dimensions. La pr diction se fait ensuite en faisant voter tous les arbres obtenus.

## 3 R sultats

Les r sultats moyens obtenus par dix r p titions de validation crois e 10-folds apparaissent dans la Table 1 compar s   des arbres classiques, des m thodes d'ensembles d'arbres et   un k-

	n	d	DT	KMDT	RF	ET	KMRF	KN
breast-cancer	683	9	0.9198	<b>0.9357</b>	0.9541	0.9537	<b>0.9546</b>	0.9601
digits	1797	64	0.8679	<b>0.9495</b>	0.9607	0.9682	<b>0.9789</b>	0.987
iris	150	4	0.9443	<b>0.9522</b>	0.945	0.9523	<b>0.9584</b>	0.9639
ringnorm	7400	20	<b>0.8976</b>	0.8252	0.954	<b>0.9638</b>	0.9056	0.6902
satellite	6435	36	0.8578	<b>0.8784</b>	0.9097	0.9083	<b>0.9099</b>	0.9087
spambase	6301	57	<b>0.9129</b>	0.7735	0.9461	<b>0.947</b>	0.9295	0.7978
vehicle	846	18	<b>0.7166</b>	0.6067	0.7495	<b>0.7603</b>	0.7134	0.6495
waveform	5000	21	0.7565	<b>0.7899</b>	0.8218	0.8239	<b>0.8362</b>	0.8197
zoo	101	16	0.8773	<b>0.908</b>	0.8993	0.9037	<b>0.9263</b>	0.782

TAB. 1 – Moyenne des scores de précision sur les ensembles de test pour différentes méthodes :  
— arbres : DT : arbre de décision (Breiman [2]), KMDT : arbre de décision k-means  
— ensembles : RF : forêt d’arbres (Breiman [3]), ET : extra-arbres (Geurts [4]), KMRF : forêt d’arbres k-means,  
— autre : KN : k-plus-proches-voisins  
et différents ensembles de données avec  $n$  le nombre d’exemples et  $d$  la dimension.

plus-proches-voisins, tous avec optimisation des paramètres. Les résultats obtenus sur certains ensembles par les KMDT dépassent de plusieurs points ceux des arbres de décision classiques.

L’un des critères déterminants pour la qualité du modèle est sa taille. Quand la taille du KMDT dépasse considérablement (4-6 fois) celle d’un DT et est de l’ordre (en nombre de noeuds total) du nombre d’exemples de la base de données, le modèle tend à être moins efficace en prédiction (ringnorm, spambase, vehicle). C’est également sur ces ensembles que la méthode des k-plus-proches-voisins est la moins efficace, reflétant une structure des données moins décisive. Dans les autres ensembles de données testés, les KMDT sont entre 1.3 et 4 fois plus grands que les DT, ce qui est inévitable puisqu’ils s’arrêtent quand ils atteignent des feuilles pures sans optimiser ce critère pour la construction, contrairement aux DT. La formation d’ensembles permet d’améliorer les performances mais moins que dans le cas des DT.

## 4 Conclusions et perspectives

La méthode permet, lorsque les classes attendues sont adaptées aux k-means (i.e "rondes") d’obtenir de bons résultats en apprentissage et un critère de taille du modèle permet de prédire son efficacité. Les résultats encourageants de ce premier travail sur des modèles de classification structurelle pour l’apprentissage automatique ouvrent la voie pour l’utilisation de modèles plus généraux comme les systèmes de classes totalement équilibrés (Lehel [6]) et vers la généralisation des arbres de décision.

## Références

- [1] Alexander Böcker, Swetlana Derksen, Elena Schmidt, Andreas Teckentrup and Gisbert Schneider. A hierarchical clustering approach for large compound libraries *J. Chem. Inf. Model.*, 45 :807–815, Janvier 2005.
- [2] Leo Breiman, Jerome H. Friedman, Richard A. Olshen, and Charles J. Stone. *Classification and Regression Trees.* Wadsworth and Brooks, Monterey, CA, 1984.
- [3] Leo Breiman. Random forests. *Mach. Learn.*, 45(1) :5–32, Octobre 2001.
- [4] Pierre Geurts, Damien Ernst, Louis Wehenkel. Extremely Randomized Trees. *Mach. Learn.*, 63(1) :3–42, Avril 2006.
- [5] Tin Kam Ho. The random subspace method for constructing decision forests. *IEEETrans. Pattern Anal. Mach. Intell.*, 20(8) :832–844,1998.
- [6] Jenő Lehel. A characterization of totally balanced hypergraphs. *Discrete Math.*, 57(1–2), 59–65, Novembre 1985.