

APPRENTISSAGE STATISTIQUE DE LA CLASSIFICATION :

PARTIE I.

G erard Biau, Laurent Rouvi re

Universit  Paris VI, Universit  Rennes 2

`biau@math.univ-montp2.fr`

`laurent.rouviere@univ-rennes2.fr`

L'*apprentissage statistique* d signe un vaste ensemble de m thodes et d'algorithmes permettant, dans un sens g n ral, d'extraire l'information pertinente de donn es ou d'apprendre des comportements   partir d'exemples. Les applications de ce paradigme sont tr s nombreuses, allant de la recherche d'informations dans de grands ensembles de donn es (fouille de textes ou d'images)   la biologie (reconstruction des r seaux g n tiques, puces ADN, etc).

Notre expos  se divisera en trois parties. Dans la premi re partie, nous rappelons les concepts d'apprentissage supervis  et non-supervis    travers divers exemples. Dans le premier cas, il s'agit de trouver une fonction g (appel e *r gle*) qui permet d'expliquer une variable   partir d'observations d'autres variables. Dans le second cas, c'est- -dire en l'absence d'une variable   expliquer, l'objectif est la recherche d'une typologie permettant de regrouper les observations en classes homog nes. Il existe de nombreuses approches permettant de traiter ces deux domaines de la statistique, comme par exemple :

- La discriminante de Fisher, les SVM (Support Vector Machines), les plus proches voisins, la r gle du noyau, les arbres de classification dans le cas supervis  ;
- La m thode des k -means dans le cas non-supervis .

Les diff rentes m thodes de classification  nonc es ci-dessus d pendent en g n ral d'un ou plusieurs param tres dont le choix se r v le crucial. Dit autrement,  tant donn e une famille de r gles candidates, comment choisir,   partir des seules observations, la "meilleure" r gle (ou encore une r gle qui ne s' loigne pas trop de l'optimum) en un sens   pr ciser ? Dans la deuxi me partie de l'expos , nous abordons cette probl matique via la th orie de Vapnik-Chervonenkis qui propose des techniques de s lection bas es sur la minimisation de crit res empiriques. Nous mettrons en particulier l'accent sur l'int r t de cette approche :

- D'un point de vue pratique, la proc dure est "automatique" : nul besoin pour le praticien de choisir certains param tres ;

- D'un point de vue théorique, il est possible d'établir des bornes de performance permettant de mesurer la qualité de la règle sélectionnée.

Dans la troisième et dernière partie de l'exposé, nous verrons comment adapter les algorithmes de l'apprentissage à des données plus complexes. En effet, dans de nombreux domaines de la statistique contemporaine, les individus prennent la forme de *courbes* (ou *surfaces*). Ces courbes peuvent, par exemple, représenter la température en un point du globe, le cours d'une action en bourse, le tracé d'un électrocardiogramme ou encore la consommation d'électricité d'une grande ville... Les appareils de mesure ne captant que la valeur de la courbe à certains instants, les données disponibles ne sont en fait que des versions discrétisées de la courbe. Toutefois, le statisticien a intérêt à prendre en compte le caractère *continu* de ces observations. Il devient dès lors pertinent de ne plus considérer ces individus comme des vecteurs de grande dimension, mais plutôt de les appréhender comme des "fonctions", c'est-à-dire comme des objets uniques évoluant dans des espaces de dimension infinie. Il faut alors adapter les méthodes statistiques classiques à ces nouveaux individus. L'approche que nous développons consiste à réduire la dimension infinie des observations en ne considérant que certains coefficients des données décomposées dans une base appropriée (Fourier ou ondelettes par exemple). Nous présentons les propriétés asymptotiques et à distance finie de la méthode envisagée et nous illustrons les performances de cette approche sur des jeux de données réelles et simulées.

Références

- Hastie, T., Tibshirani, T. et Friedman, J. (2001) *Elements of Statistical Learning*, Springer.
- Vapnik, V.N., (1995) *The Nature of Statistical Learning Theory*, Springer.