

---

# FONDEMENTS MATHÉMATIQUES DE L'APPRENTISSAGE STATISTIQUE

*par*

Christophe Giraud

---

**Résumé.** L'objectif d'un algorithme de classification est de prédire au mieux la classe d'un objet à partir d'observations de cet objet. Un exemple typique est le filtre à spam des messageries électroniques qui prédisent (plus ou moins bien) si un courriel est un spam ou non. Nous introduisons dans ces notes les principaux concepts fondamentaux de la théorie de la classification statistique supervisée et quelques uns des algorithmes de classification les plus populaires. Nous soulignons chemin faisant l'importance de certains concepts mathématiques, parmi lesquels la symétrisation, la convexification, les inégalités de concentration, le principe de contraction et les espaces de Hilbert à noyau reproduisant.

## 1. Introduction

En ce début de siècle, nous observons un accroissement phénoménal de l'utilisation des mathématiques, à la fois dans l'industrie et dans les laboratoires scientifiques. Cet essor de l'importance des mathématiques va de pair avec l'explosion des volumes de données collectées et l'accroissement de la puissance de calculs des ordinateurs. Dans l'industrie, la modélisation mathématique apparaît à tous les stades de la vie d'un produit. Depuis la conception technique, avec force de simulations numériques, via la production, avec l'optimisation des ressources et des flux, jusqu'au marketing et la distribution avec des prévisions basées sur l'analyse de grandes bases de données. Dans les laboratoires scientifiques, la modélisation mathématique devient de plus en plus cruciale, en particulier en biologie et médecine où les scientifiques doivent extraire des informations pertinentes des

données massives qu'ils produisent grâce aux récents développements biotechnologiques.

La classification automatique est peut-être l'un des usages quotidiens les plus invasifs des mathématiques. L'objectif de la classification automatique est de prédire au mieux la classe  $y$  d'un objet  $x$  à partir d'observations de ce dernier. Un exemple typique est le filtre à spam de notre messagerie électronique qui prédit (plus ou moins bien) si un courriel est un spam ou non. La classification automatique est omniprésente dans notre quotidien, filtrant nos courriels, lisant automatiquement les codes postaux sur nos lettres ou reconnaissant les visages sur les photos que nous publions sur les réseaux sociaux. Elle est aussi extrêmement importante en sciences, par exemple en médecine pour effectuer des diagnostics précoces à partir de données à hauts débits, ou pour la recherche *in silico* de médicaments efficaces.

Nous introduisons dans ces notes les principaux concepts fondamentaux de l'apprentissage statistique (supervisé). Nous décrivons dans la partie 2 la modélisation mathématique d'un problème générique de classification. Dans la partie 3, nous analysons la précision prédictive d'un algorithme universel de classification et dans la partie 4 nous dérivons de cet algorithme théorique des algorithmes numériquement implémentables et populaires. Les appendices rassemblent des résultats techniques nécessaires pour la définition et l'analyse des algorithmes de classification.

## 2. Modélisation mathématique

Par soucis de simplicité, nous allons nous restreindre au cas où il y a seulement deux classes (comme pour le filtre à spam). Le problème de la classification automatique peut alors être modélisé de la manière suivante. Soit  $\mathcal{X}$  un espace mesuré. On observe conjointement un point  $X \in \mathcal{X}$  et une étiquette  $Y \in \{-1, +1\}$ . Notre objectif est de construire une fonction  $h : \mathcal{X} \rightarrow \{-1, +1\}$ , appelée *classifieur*, telle que  $h(X)$  prédit au mieux l'étiquette  $Y$ .

Supposons que le couple  $(X, Y) \in \mathcal{X} \times \{-1, +1\}$  est issu d'un tirage selon une loi  $\mathbb{P}$ . Pour un classifieur  $h : \mathcal{X} \rightarrow \{-1, +1\}$  donné, la probabilité de mauvaise classification est

$$L(h) = \mathbb{P}(Y \neq h(X)).$$



























sont obtenus en résolvant (15) pour des choix spécifiques de  $\mathcal{F}$  et  $\ell$ , voir par exemple les Parties 4.3 et 4.4 pour quelques exemples.

*Quelques fonctions de perte convexes  $\ell$  classiques.* Il est naturel de considérer une fonction de perte  $\ell$  qui est décroissante et positive. Habituellement, on demande aussi que  $\ell(z) \geq \mathbf{1}_{z < 0}$  pour tout  $z \in \mathbb{R}$  car cela nous permet de donner une majoration de la probabilité de mauvaise classification, voir le Théorème 3. Quelques fonctions de pertes classiques sont

- la perte exponentielle  $\ell(z) = e^{-z}$
- la perte logit  $\ell(z) = \log_2(1 + e^{-z})$
- la perte hinge  $\ell(z) = (1 - z)_+$  (avec  $(x)_+ = \max(0, x)$ )

voir la Figure 2 pour un tracé de ces trois fonctions.

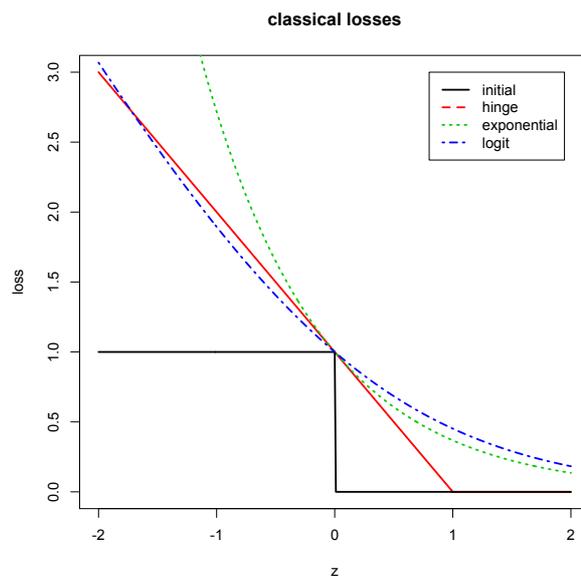


FIGURE 2. Tracé des pertes exponentielle, hinge et logit

*Quelques exemples classiques d'ensembles fonctionnels  $\mathcal{F}$ .* Les principaux exemples classiques d'ensembles fonctionnels  $\mathcal{F}$  peuvent être regroupés en deux classes.



















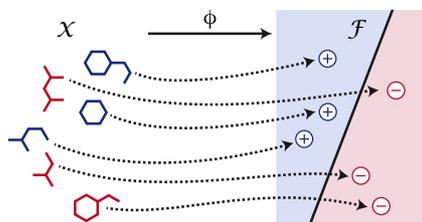


FIGURE 5. Classification de molécules à l'aide d'un SVM.

comme représenté dans la Figure 4. Cela permet d'obtenir un algorithme non linéaire pour un coût similaire à celui d'un algorithme linéaire.

La seconde raison pour utiliser un RKHS est de pouvoir employer sur n'importe quel ensemble  $\mathcal{X}$  des algorithmes définis pour des vecteurs. Imaginons par exemple qu'on veuille classifier des protéines ou des molécules en fonction de leur propriétés thérapeutiques. Notons par  $\mathcal{X}$  notre ensemble de molécules. Pour tout  $x, y \in \mathcal{X}$ , représentons par  $k(x, y)$  leur similarité (selon des critères à définir). Si le noyau résultant  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  est défini positif, on peut alors directement appliquer un SVM pour les classifier, comme illustré Figure 5. Bien sûr, le point clef dans ce cas est de construire le noyau  $k$ . En général, le noyau  $k(x, y)$  est défini en fonction de certaines propriétés de  $x, y$  qui sont connues pour avoir de l'importance pour le problème de classification. Par exemple, le nombre de courtes séquences communes est un indice utile pour quantifier la similarité entre deux protéines. La complexité du calcul de  $k(x, y)$  est aussi un critère crucial pour les applications sur des données complexes. Nous renvoyons aux excellents slides de Jean-Philippe Vert pour la description d'applications prometteuses en biologie et médecine :

<http://cbio.ensmp.fr/~jvert/talks/120302ensae/ensae>

#### 4.4. AdaBoost

AdaBoost est un algorithme qui tend à calculer une solution approximative de l'estimateur (15) avec la perte exponentielle  $\ell(z) = e^{-z}$  et l'espace fonctionnel  $\mathcal{F} = \text{span}\{h_1, \dots, h_p\}$  où  $h_1, \dots, h_p$  sont  $p$  classifieurs donnés.

Le principe de l'algorithme AdaBoost est de réaliser une minimisation dite agressive de (15). Plus précisément, AdaBoost produit une suite de fonctions  $\hat{f}_m$  pour  $m = 0, \dots, M$  en partant de  $\hat{f}_0 = 0$  puis en résolvant récursivement pour  $m = 1, \dots, M$

$$\hat{f}_m = \hat{f}_{m-1} + \beta_m h_{j_m}$$

$$\text{où } (\beta_m, j_m) = \underset{\substack{j=1, \dots, p \\ \beta \in \mathbb{R}}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \exp(-Y_i(\hat{f}_{m-1}(X_i) + \beta h_j(X_i))).$$

La classification finale est effectuée à l'aide de  $\hat{h}_M(x) = \operatorname{sign}(\hat{f}_M(x))$  qui est une approximation de  $\hat{h}_{\mathcal{H}}$  défini par (15).

La perte exponentielle permet de calculer  $(\beta_m, j_m)$  très efficacement. En effet, en posant  $w_i^{(m)} = n^{-1} \exp(-Y_i \hat{f}_{m-1}(X_i))$ , on peut écrire

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \exp(-Y_i(\hat{f}_{m-1}(X_i) + \beta h_j(X_i))) \\ = (e^\beta - e^{-\beta}) \sum_{i=1}^n w_i^{(m)} \mathbf{1}_{h_j(X_i) \neq Y_i} + e^{-\beta} \sum_{i=1}^n w_i^{(m)}. \end{aligned}$$

Lorsque la condition suivante

$$\operatorname{err}_m(j) = \frac{\sum_{i=1}^n w_i^{(m)} \mathbf{1}_{h_j(X_i) \neq Y_i}}{\sum_{i=1}^n w_i^{(m)}} \leq \frac{1}{2} \quad \text{pour tout } j = 1, \dots, p,$$

est satisfaite, les minimiseurs  $(\beta_m, j_m)$  sont donnés par

$$j_m = \underset{j=1, \dots, p}{\operatorname{argmin}} \operatorname{err}_m(j) \quad \text{et} \quad \beta_m = \frac{1}{2} \log \left( \frac{1 - \operatorname{err}_m(j_m)}{\operatorname{err}_m(j_m)} \right).$$

En remarquant que  $-Y_i h(X_i) = 2\mathbf{1}_{Y_i \neq h(X_i)} - 1$  on obtient la formulation standard de l'algorithme AdaBoost.











nous avons

$$\mathbb{E}_\sigma \left[ \sup_{z \in \mathcal{X}} \left| \sum_{i=1}^n \sigma_i \varphi(z_i) \right| \right] \leq \alpha \mathbb{E}_\sigma \left[ \sup_{z \in \mathcal{X}} \left| \sum_{i=1}^n \sigma_i z_i \right| \right].$$

Un preuve concise est donnée dans le Chapitre 11 du livre de Boucheron, Lugosi et Massart [2].

### Références

- [1] Boucheron, S., Bousquet, O. and Lugosi, G. *Theory of classification : some recent advances*. ESAIM Probability & Statistics, **9** (2005) : 323–375.
- [2] Boucheron, S., Lugosi, G. and Massart, P. *Concentration Inequalities*. Oxford University Press, 2013.
- [3] Cléménçon, S., Lugosi, G. and Vayatis, N. *Ranking and empirical risk minimization of U-statistics*. The Annals of Statistics, **36** (2008) : 844–874.
- [4] Devroye, L., Györfi, L. and Lugosi, G. *A Probabilistic Theory of Pattern Recognition*. Springer, 1996.
- [5] Hastie, T., Tibshirani, R. and Friedman, J. *The element of statistical learning*. Springer, 2009.
- [6] McDiarmid, C. *On the Method of Bounded Differences*. Surveys in Combinatorics **141** (1989) : 148–188.

---

CHRISTOPHE GIRAUD, Département de Mathématiques, Bât. 425, Faculté des Sciences d’Orsay, Université Paris-Sud, F-91405 Orsay Cedex  
 CMAP, UMR CNRS 7641, Ecole Polytechnique, Route de Saclay, 91128 Palaiseau Cedex • *E-mail* : christophe.giraud@math.u-psud.fr  
*Url* : <http://www.cmap.polytechnique.fr/~giraud/>

