

M

Massih-Reza AMINI

Préface de Francis Bach

Machine Learning

Programmes libres (GPLv3)
essentiels au développement
de solutions big data

2^e édition

Machine Learning

Machine Learning et intelligence artificielle

Le Machine Learning est l'un des domaines de l'intelligence artificielle qui a pour but de concevoir des programmes qui ne sont pas explicitement codés pour s'acquitter d'une tâche particulière. Les concepts de ce domaine sont fondés sur la logique inférentielle et tentent de dégager des règles générales à partir d'un nombre fini d'observations.

Un ouvrage de référence

Cet ouvrage présente les fondements scientifiques de la théorie de l'apprentissage supervisé, les algorithmes les plus répandus développés suivant ce domaine ainsi que les deux cadres de l'apprentissage semi-supervisé et de l'ordonnancement, à un niveau accessible aux étudiants de master et aux élèves ingénieurs. La première édition, connue sous le nom *Apprentissage machine*, fut traduite en chinois par les éditions iTuring. Dans cette deuxième édition, un nouveau chapitre est dédié au *Deep Learning*, sur les réseaux de neurones artificiels, et nous avons réorganisé les autres chapitres pour un exposé cohérent reliant la théorie aux algorithmes développés dans cette sphère. Vous trouverez également dans cette édition quelques programmes des algorithmes classiques, écrits en langages Python et C (langages à la fois simples et populaires), et à destination des lecteurs qui souhaitent connaître le fonctionnement de ces modèles désignés parfois comme des boîtes noires. Ces programmes libres (GPLv3) essentiels au développement de solutions big data sont déposés progressivement sur ce gitlab (<https://gricad-gitlab.univ-grenoble-alpes.fr/aminima/machine-learning-tools>).

À qui s'adresse ce livre ?

- Aux élèves ingénieurs, étudiants de master et doctorants en mathématiques appliquées, algorithmique, recherche opérationnelle, gestion de production, aide à la décision.
- Aux ingénieurs, enseignants-chercheurs, informaticiens, industriels, économistes et décideurs ayant à résoudre des problèmes de classification, de partitionnement et d'ordonnancement à large échelle.

Massih-Reza Amini, professeur d'informatique à l'université Grenoble Alpes (UGA), est titulaire d'une thèse sur l'étude de nouveaux cadres et modèles d'apprentissage statistiques pour les nouvelles applications émergentes issues d'Internet.

Il est co-auteur de plus d'une centaine d'articles scientifiques parus parmi les actes de conférences et de revues les plus prestigieux des domaines de l'apprentissage automatique et de la recherche d'information. Il est également co-auteur des ouvrages Recherche d'information et Data science parus aux éditions Eyrolles.

Éditions
EYROLLES

www.editions-eyrolles.com
Éditions Eyrolles | Diffusion Geodif

Code éditeur : G67947
ISBN : 978-2-212-67947-2

Machine Learning

**Programmes libres (GPLv3)
essentiels au développement
de solutions big data**

SUR LE MÊME THÈME

P. BOSC, M. GUYOMARD, L. MICLET. – **Conception d’algorithmes.**
N°67728, 2^e édition, 2019, 852 pages.

M.-R. AMINI, R. BLANCH, M. CLAUSEL, J.-B. DURAND, E. GAUSSIER, J. MALICK,
C. PICARD, C. QUEMA, G. QUENOT. – **Data science : cours et exercices.**
N°67410, 2018, 272 pages.

A. CORNUEJOLS, L. MICLET. – **Apprentissage artificiel.**
N°67522, 3^e édition, 2018, 912 pages.

E. BIERNAT, M. LUTZ. – **Data science : fondamentaux et études de cas.**
N°14243, 2015, 296 pages.

W. MCKINNEY. – **Analyse de données en Python.**
N°14109, 2015, 488 pages.

M.-R. AMINI, E. GAUSSIER. – **Recherche d’information.**
N°13532, 2013, 234 pages.

R. BRUCHEZ. – **Les bases de données NoSQL et le Big Data.**
N°14155, 2015, 322 pages.

Retrouvez nos bundles (livres papier + e-book) et livres numériques sur
<http://izibook.eyrolles.com>

Massih-Reza Amini

Préface de Francis Bach

Machine Learning

**Programmes libres (GPLv3)
essentiels au développement
de solutions big data**

2^e édition

● **Éditions
EYROLLES**

ÉDITIONS EYROLLES
61, bd Saint-Germain
75240 Paris Cedex 05
www.editions-eyrolles.com

En application de la loi du 11 mars 1957, il est interdit de reproduire intégralement ou partiellement le présent ouvrage, sur quelque support que ce soit, sans l'autorisation de l'Éditeur ou du Centre Français d'exploitation du droit de copie, 20, rue des Grands Augustins, 75006 Paris.

© Groupe Eyrolles, 2015,

© Éditions Eyrolles, 2020, ISBN : 978-2-212-67947-2

Préface

Depuis quelques années, dans les domaines scientifiques, industriels et personnels, la présence de données numériques et leur utilisation ont explosé. Certaines de ces données sont massives, nécessitent des outils et architectures spécifiques, comme en astronomie ou pour les moteurs de recherche, et constituent les problèmes dits de «big data».

D'autres données ne sont pas si massives, comme les photos ou vidéos familiales, mais constituent toujours un défi algorithmique. Le grand changement récent est non seulement la taille, mais aussi le côté omniprésent de ces données, qui sont utilisées quotidiennement.

Depuis une vingtaine d'années, l'apprentissage statistique («machine learning» en anglais) s'est considérablement développé, à l'interface entre l'informatique et les statistiques, et constitue le cœur méthodologique des algorithmes modernes de traitement de données. Même si les recherches en apprentissage sont toujours en plein essor, un socle méthodologique et algorithmique a émergé.

Ce livre constitue une introduction équilibrée aux concepts et outils les plus importants de l'apprentissage supervisé et de ses extensions. Un accent remarquable est mis sur des résultats théoriques élégants, simples mais puissants, des algorithmes efficaces qui ont montré leurs preuves en pratique, et des codes permettant de reproduire les expériences.

Francis Bach

octobre 2014

Table des matières

Table des figures	xii
Liste des algorithmes	xvii
Avant-propos	1
Concepts étudiés	1
Organisation du livre	3
CHAPITRE 1	
Principes de base en apprentissage supervisé	7
1.1 Principe de la minimisation du risque Empirique	9
1.1.1 Hypothèse et définitions	9
1.1.2 Énoncé du principe	11
1.2 Consistance du principe MRE	12
1.2.1 Définition	12
1.2.2 Étude de pire cas	13
1.3 Principe de la Minimisation du Risque Structurel	14
1.3.1 Estimation de l'erreur de généralisation sur un ensemble de test ..	14
1.3.2 Borne uniforme sur l'erreur de généralisation	16
1.3.3 Énoncé du principe	27
CHAPITRE 2	
Bornes de généralisation dépendantes des données	29
2.1 Complexité de Rademacher	30
2.2 Lien entre la complexité de Rademacher et la dimension VC	31
2.2.1 Différentes étapes d'obtention d'une borne de généralisation avec la complexité de Rademacher	33
2.2.2 Propriétés de la complexité de Rademacher	38

2.3	Considérations pratiques	41
2.3.1	Régularisation	41
2.3.2	Minimisation d'une borne convexe de l'erreur de classification	42
2.4	Cas multi-classe	44
2.4.1	Erreurs de classification	45
2.4.2	Réduction du problème multi-classe à la classification binaire	46
2.4.3	Borne sur l'erreur de généralisation	50
CHAPITRE 3		
	Algorithmes d'optimisation à direction de descente	55
3.1	Algorithme du gradient	57
3.1.1	Mode batch	58
3.1.2	Mode en ligne	60
3.2	Méthode de quasi-Newton	61
3.2.1	Direction de Newton	62
3.2.2	Formule de Broyden-Fletcher-Goldfarb-Shanno	63
3.3	Recherche linéaire	66
3.3.1	Conditions de Wolfe	67
3.3.2	Algorithme de recherche linéaire basé sur une stratégie de retour en arrière	72
3.4	Méthode du gradient conjugué	74
3.4.1	Directions conjuguées	74
3.4.2	Algorithme du gradient conjugué	77
CHAPITRE 4		
	Deep Learning	81
4.1	Perceptron	84
4.1.1	Théorème de convergence du perceptron	89
4.1.2	Perceptron à marge et lien avec le principe MRE	91
4.2	Adaline	92
4.2.1	Lien avec la régression linéaire et le principe MRE	92
4.2.2	Différence avec le perceptron	94
4.3	Régression logistique	95
4.3.1	Lien avec le principe MRE	95
4.3.2	Modèle formel	96
4.4	Perceptron multi-couche	97
4.4.1	Représentation formelle	97
4.4.2	Algorithme de rétropropagation de l'erreur	99
4.5	Réseaux convolutifs et récurrents	103
4.5.1	Réseaux convolutifs	103

4.5.2 Réseaux récurrents	106
4.6 Considérations pratiques	108
4.6.1 Fonctions de transfert	108
4.6.2 Méthode du moment	109
4.6.3 Traitement par mini-batches	109
4.6.4 Normalisation par batchs	110
4.6.5 Technique de décrochage	111
CHAPITRE 5	
Séparateurs à Vaste Marge	113
5.1 Notion de marge	114
5.1.1 Marge dure	114
5.1.2 Marge souple	119
5.2 Astuce du noyau	127
5.2.1 Définition de fonction noyau	127
5.2.2 Noyau symétrique défini positif	128
5.2.3 SVM avec des noyaux symétriques définis positifs	129
5.3 Étude théorique et cas multi-classe	129
5.3.1 Borne de généralisation à base de marge	129
5.3.2 Séparateurs à vaste marge multi-classe	132
CHAPITRE 6	
Boosting	139
6.1 Adaboost	140
6.1.1 Lien avec le principe MRE	143
6.1.2 Échantillonnage par rejet	144
6.2 Étude théorique	145
6.2.1 Borne sur l'erreur empirique à base de marge	146
6.2.2 Borne de généralisation à base de marge du classifieur de vote	148
6.3 AdaBoost multi-classe	150
6.3.1 Pseudo-erreur de classification	150
6.3.2 Échantillonnage par rejet suivant deux distributions	152
CHAPITRE 7	
Apprentissage semi-supervisé	155
7.1 Cadre non supervisé et hypothèses de base	156
7.1.1 Mélange de densités	156
7.1.2 Estimer les paramètres du mélange	157
7.1.3 Hypothèses de base en apprentissage semi-supervisé	165
7.2 Méthodes génératives	167

7.2.1	Extension des critères à base de vraisemblance au cas semi-supervisé	168
7.2.2	Algorithme CEM semi-supervisé	169
7.2.3	Application : apprentissage semi-supervisé d'un classifieur Naive Bayes	170
7.3	Méthodes discriminantes	173
7.3.1	Algorithme auto-apprenant	174
7.3.2	Séparateurs à vaste marge transductifs	176
7.3.3	Borne transductive sur l'erreur du classifieur de Bayes	179
7.3.4	Apprentissage multi-vue basé sur le pseudo-étiquetage	183
7.4	Méthodes graphiques	186
7.4.1	Propagation des étiquettes	186
7.4.2	Marche aléatoire markovienne	189
CHAPITRE 8		
	Apprentissage de fonctions d'ordonnement	193
8.1	Formalisme	194
8.1.1	Fonctions d'erreur d'ordonnement	194
8.1.2	Ordonnement d'instances	198
8.1.3	Ordonnement d'alternatives	200
8.2	Approches	203
8.2.1	Par point	203
8.2.2	Par paire	208
8.3	Apprentissage avec des données interdépendantes	217
8.3.1	Borne de test	219
8.3.2	Borne de généralisation	220
8.3.3	Estimation des bornes pour quelques exemples d'application	226
ANNEXE A		
	Rappels de probabilités	235
A.1	Mesure de probabilité	235
A.1.1	Espace probabilisable	235
A.1.2	Espace probabilisé	236
A.2	Probabilité conditionnelle	237
A.2.1	Formule de Bayes	237
A.2.2	Indépendance en probabilité	239
A.3	Variables aléatoires réelles	239
A.3.1	Fonction de répartition	240
A.3.2	Espérance et variance d'une variable aléatoire	241
A.3.3	Inégalités de concentration	242

ANNEXE B	
Codes programmes	247
B.1 Structures de données	247
B.1.1 Base de données	247
B.1.2 Structure des hyper-paramètres	248
B.2 Structure pour une représentation creuse	249
B.3 Lancement des programmes	251
B.4 Codes	253
B.4.1 Algorithme BGFS (chapitre 3, section 3.2.2)	253
B.4.2 Recherche linéaire (chapitre 3, section 3.3)	256
B.4.3 Gradient conjugué (chapitre 3, section 3.4)	258
B.4.4 Perceptron (chapitre 4, section 4.1)	260
B.4.5 Adaline (chapitre 4, section 4.2)	261
B.4.6 Régression logistique (chapitre 4, section 4.3)	262
B.4.9 Perceptron multi-couche (chapitre 4, section 4.4)	264
B.4.7 AdaBoost (chapitre 6, section 6.1)	267
B.4.8 AdaBoost M2 (chapitre 6, section 6.3)	270
B.4.10 K-moyennes (chapitre 7, section 7.1.2)	273
B.4.11 Naïve-Bayes semi-supervisé (chapitre 7, section 7.2.3)	275
B.4.12 Auto-apprentissage (chapitre 7, section 7.3.1)	278
B.4.13 Auto-apprentissage à une passe (chapitre 7, section 7.3.1)	281
B.4.14 PRank (chapitre 8, section 8.2.1)	282
B.4.15 RankBoost (ordonnancement bipartite - chapitre 8, section 8.2.2)	284
Bibliographie	287
Index	301

Table des figures

I	Illustration des deux phases d'un problème d'apprentissage supervisé. Dans la phase d'apprentissage (schématisée par les traits pleins), une fonction minimisant l'erreur empirique sur une base d'entraînement est trouvée parmi une classe de fonctions prédéfinies. Dans la phase de test (schématisée par les traits pointillés), les sorties de nouveaux exemples sont prédites par la fonction de prédiction.	2
1.1	Description schématique de la notion de consistence. L'axe des abscisses représente la classe des fonctions \mathcal{F} et les courbes d'erreurs empirique (en pointillé) et de généralisation (en trait plein) sont montrées en fonction de $f \in \mathcal{F}$. Le principe MRE consiste à trouver la fonction f_S dans la classe \mathcal{F} qui minimise l'erreur empirique sur une base d'entraînement S . Ce principe est consistant si en probabilité $\hat{\mathcal{L}}(f_S, S)$ converge vers $\mathfrak{L}(f_S)$ et $\inf_{g \in \mathcal{F}} \mathfrak{L}(g)$	13
1.2	Pulvérisation des points dans le plan de dimension $d = 2$ par une classe de fonctions linéaires. Chaque classifieur linéaire sépare le plan en deux sous-espaces, avec un vecteur normal qui pointe vers le sous-espace contenant les exemples appartenant à la classe +1 (représentés par des cercles pleins). Le nombre maximal de points dans le plan pouvant être pulvérisés par la classe de fonctions linéaires, ou la dimension VC de cette classe de fonctions, est dans ce cas égal à 3.	22
1.3	Construction des ensembles \mathcal{F}_1 et \mathcal{F}_2 à partir de la classe de fonctions \mathcal{F} pour la preuve du lemme de Sauer (1972) sur un exemple jouet.	23
1.4	Illustration du principe de la minimisation du risque structurel. L'axe des abscisses montre une hiérarchie de sous-ensembles de fonctions imbriqués avec une capacité croissante, de gauche à droite. Plus la capacité d'une classe de fonctions est grande, plus l'erreur empirique d'une fonction de cette classe sur une base d'entraînement sera faible et plus la borne sur son erreur de généralisation sera mauvaise. Le principe de la minimisation du risque structurel consiste à choisir la fonction de la classe de fonctions pour laquelle nous avons la meilleure estimation de sa borne de généralisation.	27

2.1	Quatre fonctions de coût pour un problème de classification à deux classes en fonction du produit $y \times h$, où h est la fonction apprise. Les fonctions de coût sont l'erreur instantanée de classification : $\mathbb{1}_{y \times h \leq 0}$, le coût exponentiel : $e^{-y \times h}$, le coût logistique : $\frac{1}{\ln(2)} \times \ln(1 + \exp(-y \times h))$ et le coût de hinge : $\max(0, 1 - y \times h)$. Des valeurs positives (négatives) de $y \times h$ impliquent une bonne (mauvaise) classification.	44
2.2	Coût à base de marge (en gris plein), majorant l'erreur de classification (en noir) et défini avec le paramètre de marge ρ	51
3.1	Illustration d'une ligne de niveau elliptique d'une fonction de coût, $\mathcal{L}(\mathbf{w})$ continue et doublement dérivable au voisinage de son minimum \mathbf{w}^* . Au voisinage de ce minimum, les axes des ellipses sont définis par rapport aux vecteurs propres de la matrice hessienne de \mathcal{L} calculés au point \mathbf{w}^* et ils sont inversement proportionnels aux racines carrées des valeurs propres associées à ces vecteurs propres.	57
3.2	Illustration de la minimisation d'une fonction d'erreur convexe $\mathcal{L}(\mathbf{w})$ avec l'algorithme du gradient (équations 3.5). Les courbes elliptiques représentent les lignes de niveau sur lesquelles la fonction d'erreur admet des valeurs constantes. Les vecteurs \mathbf{v}_1 et \mathbf{v}_2 représentent les vecteurs propres de la hessienne (le minimum recherché est au centre des axes). On remarque le mouvement d'oscillation des poids trouvés autour du minimum, et on note que l'opposé du gradient de la fonction de coût estimé à ces points ne pointe généralement pas vers ce même minimum.	58
3.3	Illustration de la direction de l'opposé du gradient, $\mathbf{p}_t = -\nabla \mathcal{L}(\mathbf{w}^{(t)})$ et celle de Newton $\mathbf{p}_t = -\mathbf{H}^{-1} \nabla \mathcal{L}(\mathbf{w}^{(t)})$ qui, à partir de n'importe quel point $\mathbf{w}^{(t)}$ sur la fonction de coût, pointe vers le minimiseur \mathbf{w}^*	62
3.4	Illustration du cas non convergent d'une séquence décroissante des poids $(\mathbf{w}^{(t)})_{t \in \mathbb{N}}$ vers le minimiseur d'une fonction objectif \mathcal{L} , lorsque la décroissance est trop faible par rapport aux longueurs des sauts.	68
3.5	Illustration du cas non convergent d'une séquence décroissante des poids $(\mathbf{w}^{(t)})_{t \in \mathbb{N}}$ vers le minimiseur d'une fonction objectif \mathcal{L} , lorsque les sauts sont trop petits par rapport au taux initial de la décroissance.	69
3.6	Les valeurs admissibles du pas d'apprentissage selon les critères de Wolfe (équations 3.16 et équations 3.17) sur un exemple jouet. La pente de la tangente à la courbe $\eta \mapsto \mathcal{L}(\mathbf{w}^{(t)} + \eta \mathbf{p}_t)$ au point $\eta = 0$ est $\mathbf{p}_t^\top \nabla \mathcal{L}(\mathbf{w}^{(t)}) < 0$, pour une valeur de $\alpha > 0$ fixée, la contrainte (équations 3.16) est vérifiée lorsque $\eta \in (0, \eta_2]$, et pour une valeur de $\beta > 0$ donnée, la contrainte de courbure (équations 3.17) est vérifiée pour $\eta \geq \eta_1$	70
3.7	Illustration du fonctionnement de l'algorithme de recherche linéaire suivant l'opposé du gradient (en gris) et la méthode du gradient conjugué (en pointillé noir).	75

4.1	Illustration d'un neurone formel de McCulloch et Pitts (1943).	84
4.2	Illustration du modèle perceptron de Rosenblatt (1958).	85
4.3	Illustration de la règle de mise à jour de l'algorithme du perceptron.	88
4.4	Illustration de l'objectif du perceptron à marge. Le vecteur normal de la frontière de décision pointe vers le demi-espace contenant des exemples positifs (en cercles pleins). Les exemples intervenant dans la modification des poids du modèle sont encerclés.	91
4.5	Modèle formel associé au perceptron et à l'adaline.	93
4.6	Illustration des solutions trouvées par les algorithmes du perceptron (en pointillé) et de l'adaline (en trait plein) pour un problème de classification linéairement séparable.	94
4.7	La différence entre l'adaline (gauche) et la régression logistique (droite). Les valeurs de sortie des deux modèles sont représentées par les surfaces, courbe pour la régression logistique et plane pour l'adaline. Les sorties prédites par le modèle logistique sont bornées par 0 et 1, alors que, pour Adaline, elles croissent, en valeur absolue, par rapport à la distance des points à la frontière de décision.	96
4.8	Architecture d'un perceptron multi-couche à n couches cachées. Sur cet exemple, les paramètres des biais sont introduits par des poids liés à deux unités supplémentaires associés à la couche d'entrée et à chaque couche cachée ayant respectivement les valeurs fixées $z_0^{(c)} = 1, \forall c = \{0, \dots, n\}$. Un réseau avec plus de deux couches cachées est appelé un réseau profond.	98
4.9	Illustration schématique des phases de propagation (en ligne continue) et de rétro-propagation (en ligne hachurée) de l'algorithme de rétro-propagation de l'erreur. Pour un exemple \mathbf{x} en entrée, la valeur de l'unité j d'une couche cachée est déterminée sur la base d'une transformation $\bar{H} : \mathbb{R} \rightarrow \mathbb{R}$ du produit scalaire entre le vecteur des valeurs des unités se trouvant sur la couche $Av(j)$ précédant la couche de l'unité j et le vecteur des poids reliant l'unité j à ces unités : $a_j = \sum_{i \in Av(j)} z_i w_{ji}$ (phase de propagation). Dans la phase de rétro-propagation, les erreurs des unités se trouvant sur la couche succédant à la couche de l'unité j , $Ap(j)$, sont combinées avec les poids reliant l'unité j à ces unités pour déterminer l'erreur δ_j associée à cette unité : $\delta_j = \bar{H}'(a_j) \sum_{l \in Ap(j)} \delta_l \times w_{lj}$	100
4.10	Un exemple de réseau convolutif pour la reconnaissance de chiffres manuscrits avec deux couches de convolution et deux couches de regroupement de c_1 et c_2 obtenues avec un noyau 5×5 et un max-pooling 2×2	106
4.11	Illustration d'un réseau récurrent.	107

5.1	Hyperplans pour un problème de classification linéairement séparable en dimension 2. Les vecteurs de support appartenant aux hyperplans marginaux d'équations $\langle \bar{\mathbf{w}}, \mathbf{x} \rangle + w_0 = \pm 1$ sont encerclés.	117
5.2	Hyperplans linéaires pour un problème de classification non linéairement séparable. Les vecteurs de support sont encerclés. Soit ces vecteurs reposent sur un des hyperplans marginaux, soit ce sont des points aberrants. La distance d'un point aberrant \mathbf{x} à l'hyperplan marginal associé à sa classe est $\frac{\xi}{\ \bar{\mathbf{w}}\ }$	120
6.1	Illustration du fonctionnement de l'algorithme d'Adaboost sur un problème jouet où la combinaison finale des apprenants faibles linéaires conduit à un classifieur non linéaire. Le vecteur normal de chaque classifieur faible pointe vers le demi-espace des exemples positifs (en cercle plein, et les exemples mal classés sont encerclés). À une itération donnée, le classifieur courant tente de bien classer les exemples mal classés par le classifieur de l'iteration précédente, et les poids \mathbf{w} de ces classifieurs, intervenant dans le vote majoritaire final, sont inversement proportionnels à leurs erreurs de classification.	142
6.2	Illustration de la technique d'échantillonnage par rejet	145
6.3	Courbe représentative de la fonction $z \mapsto (1 + 2z)^{1+z}(1 - 2z)^{1-z}$ sur l'intervalle $]0, \frac{1}{2}[$	148
7.1	Illustration des deux étapes de l'estimation et de la maximisation de l'algorithme EM. L'axe des abscisses représente les valeurs possibles de l'ensemble des paramètres Θ et l'axe des ordonnées donne le logarithme de la vraisemblance des données. L'algorithme EM calcule la fonction $Q(\Theta, \Theta^{(t)})$ en utilisant l'estimation courante $\Theta^{(t)}$ et donne la nouvelle $\Theta^{(t+1)}$ comme le point maximum de $Q(\Theta, \Theta^{(t)})$	159
7.2	Illustration des hypothèses de partition et de variété. Les exemples non étiquetés sont représentés par des cercles vides et des exemples étiquetés des différentes classes par un carré ou un triangle plein. Les partitions d'exemples sont considérées comme des régions à haute densité et la frontière de décision devrait ainsi passer par des régions de basse densité (hypothèse de partition, figure (a)). Pour un problème donné, les exemples sont supposés se trouver sur une variété géométrique de dimension inférieure (une variété de dimension 2, ou une surface courbe sur l'exemple donné - hypothèse de variété, figure (b)).	167

- 7.3 Illustration schématique de l’algorithme auto-apprenant avec un critère de marge (Tür *et al.* 2005). Un classifieur $h : \mathcal{X} \rightarrow \mathbb{R}$ est d’abord appris sur une base d’entraînement étiquetée S . L’algorithme assigne ensuite des pseudo-étiquettes de classes aux exemples non étiquetés de l’ensemble $X_{\mathcal{U}}$ en seillant les prédictions du classifieur h sur ces exemples (en gris). Les exemples pseudo-étiquetés sont retirés de l’ensemble $X_{\mathcal{U}}$ et ajoutés à l’ensemble $\tilde{S}_{\mathcal{U}}$ et un nouveau classifieur est appris en utilisant les deux ensembles S et $\tilde{S}_{\mathcal{U}}$. Ces procédés d’apprentissage et de pseudo-étiquetage sont répétés jusqu’à convergence. 176
- 7.4 Illustration des hyperplans trouvés par les algorithmes SVM (à gauche) et SVMT (à droite), pour un problème de classification binaire, les exemples étiquetés sont représentés par des cercles et des carrés et les exemples non étiquetés par des astérisques. L’algorithme SVM trouve l’hyperplan séparateur, en ignorant les exemples non étiquetés, alors que le SVMT trouve la solution qui sépare les deux classes en ne traversant pas par les régions denses. 178
- 8.1 Différentes mesures d’erreur calculées pour une fonction de score f et une requête donnée q sur un exemple jouet, ainsi que la courbe *ROC* correspondante (droite). τ_k est le taux de documents non pertinents ordonnés avant le rang k . La Précision moyenne est dans ce cas, $\mathbf{e}_{pm}(h(\mathfrak{S}, q), \mathbf{y}) = \frac{1}{4}(1 + \frac{1}{2} + \frac{3}{5} + \frac{1}{2}) = \frac{13}{20}$, et l’aire sous la courbe *ROC* vaut $\mathbf{e}_{auc}(h(\mathfrak{S}, q), \mathbf{y}) = 1 - \frac{8}{4 \times 6} = \frac{2}{3}$. Nous remarquons que c’est le rang des exemples pertinents dans la liste ordonnée induite par les scores, et non pas leurs valeurs prédites, qui intervient dans le calcul de ces mesures. 197
- 8.2 Illustration du cadre de l’ordonnancement d’alternatives. Dans la phase d’apprentissage, une fonction de score h est apprise sur un ensemble d’entrées $\{(q_1, \mathbf{y}_1), \dots, (q_m, \mathbf{y}_m)\}$ où à chaque entrée q_i est associée une liste d’alternatives $(x_1^{(i)}, \dots, x_{m_i}^{(i)})$ et des jugements de pertinence correspondants $\mathbf{y}_i = (y_1^{(i)}, \dots, y_{m_i}^{(i)})$. Une fois la fonction h trouvée, pour une nouvelle entrée q_t , sa liste d’alternatives est ordonnée suivant les valeurs de scores assignés par $h(\cdot)$ 201
- 8.3 Exemple de trois recouvrements d’un ensemble transformé $\mathfrak{T}(S)$ d’exemples interdépendants correspondant à un problème d’ordonnancement bipartite. En (a) les trois sous-ensembles \mathfrak{M}_1 , \mathfrak{M}_2 et \mathfrak{M}_3 forment un recouvrement de $\mathfrak{T}(S)$. En (b) les ensembles $\{\mathfrak{M}_j, \omega_j\}_{j \in \{1,2,3\}}$ forment un recouvrement fractionnaire de $\mathfrak{T}(S)$ et en (c) les ensembles $\{\mathfrak{M}_j, \omega_j\}_{j \in \{1,2,3\}}$ forment un recouvrement propre exact de $\mathfrak{T}(S)$ 218

Liste des algorithmes

1	Principe de la minimisation du risque structurel	28
2	Validation croisée à K plis	42
3	Stratégie un contre tous pour la classification multi-classe	47
4	Stratégie un contre un pour la classification multi-classe	48
5	Stratégie codes correcteurs d'erreur pour la classification multi-classe	49
6	Gradient stochastique	61
7	Quasi-Newton	66
8	Recherche linéaire	73
9	Gradient conjugué	78
10	Perceptron	88
11	Perceptron multi-couche	101
12	SVM à marge dure (formulation duale)	119
13	SVM à marge souple (formulation duale)	122
14	Pegasos Shalev-Shwartz <i>et al.</i> (2011)	125
15	SVM multi-classe	137
16	Adaboost	141
17	Échantillonnage par rejet	145
18	Technique d'échantillonnage par rejet suivant deux distributions	151
19	Adaboost multi-classe M2	153
20	EM	159
21	CEM	161
22	K-moyennes	163
23	CEM semi-supervisé	169
24	Auto-apprentissage multi-vue	185
25	Propagation des étiquettes pour l'apprentissage semi-supervisé	188

26	PRank	229
27	RankBoost bipartite	230
28	Recherche des règles de base	231
29	Adaptation de RankBoost à l'ordonnancement d'alternatives	232

Avant-propos

L'apprentissage machine est l'un des domaines phares de l'intelligence artificielle. Il concerne l'étude et le développement de modèles quantitatifs permettant à un ordinateur d'accomplir des tâches sans qu'il soit explicitement programmé à les faire. Apprendre dans ce contexte revient à reconnaître des formes complexes et à prendre des décisions intelligentes. Compte tenu de toutes les entrées existantes, la difficulté d'accomplir cette tâche réside dans le fait que l'ensemble des décisions possibles est généralement très complexe à énumérer. Pour contourner cette difficulté, les algorithmes en apprentissage machine ont été conçus dans le but d'acquérir de la connaissance sur le problème à traiter en se basant sur un ensemble de données limitées issues de ce problème.

Concepts étudiés

Pour illustrer ce principe, considérons le cadre de l'apprentissage supervisé que nous allons en partie traiter dans cet ouvrage. Suivant ce cadre, la décision à prendre sur une entrée donnée est prise d'après la sortie d'une fonction de prédiction qui est inférée en utilisant un ensemble d'exemples étiquetés (ou base d'entraînement), où chacun de ces exemples est une paire constituée du vecteur représentatif d'une observation dans un espace vectoriel donné, et d'une réponse associée à l'exemple (aussi appelée sortie désirée ou sortie réelle). Après la phase d'estimation ou d'apprentissage, la fonction renvoyée par l'algorithme doit permettre de prédire la réponse associée à de nouvelles observations. L'hypothèse sous-jacente dans ce cas est que les exemples sont, d'une façon générale, représentatifs du problème de prédiction sur lequel la fonction sera appliquée. En pratique, une fonction d'erreur mesure l'écart entre la prédiction du modèle sur un exemple et sa sortie désirée. À partir d'un ensemble d'entraînement donné, l'algorithme d'apprentissage choisit alors une fonction, issue d'un ensemble de fonctions défini au préalable, qui réalise l'erreur moyenne la plus faible sur les exemples de la base d'entraînement.

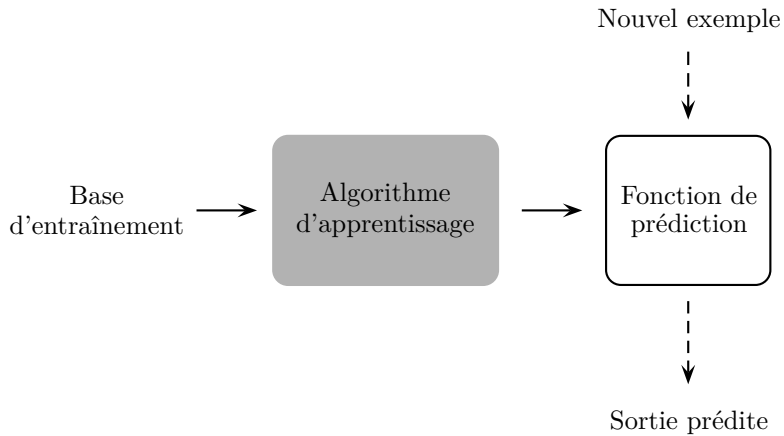


Figure I - Illustration des deux phases d'un problème d'apprentissage supervisé. Dans la phase d'apprentissage (schématisée par les traits pleins), une fonction minimisant l'erreur empirique sur une base d'entraînement est trouvée parmi une classe de fonctions prédéfinies. Dans la phase de test (schématisée par les traits pointillés), les sorties de nouveaux exemples sont prédites par la fonction de prédiction.

Cette erreur n'est généralement pas représentative de la performance de l'algorithme sur de nouveaux exemples. Il est alors nécessaire de disposer d'un second ensemble d'exemples étiquetés, ou une base de test, auxquels l'algorithme n'avait pas accès et d'estimer l'erreur moyenne de la fonction produite lors de la phase d'estimation qui sera cette fois représentative de son erreur de généralisation. On attend de l'algorithme d'apprentissage qu'il trouve une fonction ayant de bonnes performances en généralisation et non celle qui sera capable de reproduire parfaitement les réponses associées aux exemples d'entraînement (voir figure I). Les garanties d'apprenabilité du procédé de la minimisation du risque empirique ont été étudiées dans la théorie de l'apprentissage machine largement initiée par [Vapnik \(1999\)](#). Elles dépendent de la taille de la base d'entraînement et de la complexité de la classe de fonctions où on cherche la fonction de prédiction.

Historiquement, les deux tâches principales du cadre de l'apprentissage supervisé étaient la classification et la régression. Ces tâches sont similaires à la différence de l'espace des sorties désirées des exemples. Dans le cas de la classification, l'espace de sortie est discret alors qu'en régression cet espace est réel.

À la fin des années 1990 et sous l'impulsion de nouvelles technologies, notamment celles liées au développement d'Internet, de nouveaux cadres d'apprentissage ont vu le jour. Un de ces cadres est l'apprentissage avec des données partiellement

étiquetées, ou l'apprentissage semi-supervisé, dont le développement est motivé par l'effort qu'il faut consentir à constituer des bases d'apprentissage étiquetées et le constat que les données étiquetées sont chères à obtenir alors que les données non étiquetées sont foison et qu'elles contiennent de l'information sur le problème que l'on cherche à résoudre. De ce constat sont nés plusieurs travaux qui avaient pour objectif d'employer une petite quantité de données étiquetées, simultanément avec une grande quantité de données non étiquetées, pour apprendre une fonction de prédiction.

L'autre cadre qui a suscité de nombreux travaux dans la communauté d'apprentissage depuis les années 2000 concerne le développement de modèles d'ordonnement. Ce cadre a formalisé dans un premier temps les problèmes de la Recherche d'Information et il a été par la suite étendu à d'autres problèmes plus généraux.

Depuis de nombreuses années, les algorithmes d'apprentissage développés suivant ces cadres ont été appliqués avec succès à une grande variété de problèmes, incluant la reconnaissance de la parole et de l'écriture manuscrite, la vision par ordinateur, la prédiction de la structure des protéines, les systèmes de recommandations, la classification documentaire, les moteurs de recherche, etc.

Organisation du livre

Cet ouvrage présente les fondements scientifiques de la théorie de l'apprentissage supervisé, les algorithmes les plus répandus développés suivant ce cadre ainsi que les deux cadres de l'apprentissage évoqués plus haut, à un niveau accessible aux étudiants de master et aux élèves ingénieurs. Notre souci a été de proposer un exposé cohérent reliant la théorie aux algorithmes développés dans ce domaine. En outre, cette étude ne se limite pas à l'exposé de ces fondements, mais présente aussi quelques programmes des algorithmes classiques proposés dans ce manuscrit, écrits dans un langage informatique simple et populaire qui est le langage C¹, et à destination des lecteurs qui cherchent à connaître le fonctionnement de ces modèles désignés parfois comme des boîtes noires. Ce livre est organisé en six chapitres principaux et deux annexes. L'enchaînement des idées présentées dans chacun d'eux est le suivant :

- Dans le chapitre 1, nous décrivons les concepts fondamentaux de la théorie de l'apprentissage statistique de Vapnik (1999). Nous exposons la notion de consistance du principe de la minimisation du risque empirique selon lequel la plupart des algorithmes en apprentissage supervisé ont été développés. L'étude de cette consistance nous mènera à l'exposé du second principe fondamental en apprentissage qui est la minimisation du risque structurel, ouvrant le champ au développement de nouveaux modèles en apprentissage machine. En particulier,

1. <http://www.tiobe.com/index.php/content/paperinfo/tpci/index.html>

nous présentons dans ce chapitre la notion de borne sur l'erreur de généralisation en décrivant les hypothèses et les outils nécessaires pour l'obtenir.

- Dans le chapitre 2, nous allons présenter des bornes sur l'erreur de généralisation qui peuvent être estimées sur une base d'entraînement qui sert à apprendre un modèle. Ces bornes sont basées sur une notion de complexité de classes dépendante des données, appelée complexité de Rademacher. Avec cette notion il est aussi possible de dériver facilement des bornes sur l'erreur de généralisation pour la classification multi-classe. Nous exposons les approches multi-classe basées sur la classification binaire, appelées approches combinées, et dérivons une borne sur l'erreur de généralisation des classifieurs linéaires dans ce cas.
- C'est dans le chapitre 3 que nous nous intéressons aux algorithmes de base issus du domaine de l'optimisation pour la minimisation d'une fonction de risque convexe majorant le risque empirique appelés *algorithmes à direction de descente*. En particulier, nous présentons les conditions nécessaires à vérifier un algorithme à direction de descente pour converger vers le minimiseur d'une fonction objectif convexe et nous décrivons quelques variantes simples et efficaces de cet algorithme.
- Le chapitre 4 présente les principaux modèles qui sont issus des travaux sur la modélisation numérique des neurones du système nerveux, appelés plus communément réseaux de neurones artificiels ou *Deep Learning*. Ces modèles étaient précurseurs au développement des méthodes quantitatives en Intelligence Artificielle et un soin particulier a été porté à la présentation de ces modèles dans leur contexte historique.
- Dans le chapitre 5, nous présentons les séparateurs à vaste marge (SVM) qui sont issus du principe de la minimisation du risque structurel. Ces modèles sont devenus très populaires grâce à leurs justifications théoriques. En particulier, nous verrons comment utiliser l'astuce de noyau pour plonger l'espace d'entrée dans un espace de plus grande dimension dans lequel le problème d'apprentissage devient plus simple à résoudre et nous présentons l'extension de ces au cas multi-classe.
- Le chapitre 6 présente l'algorithme Adaboost issus des travaux de [Schapire \(1999\)](#). Cet algorithme combine plusieurs classifieurs de base, appelés apprenants faibles, pour construire un classifieur final, appelé apprenant fort, qui est plus performant que chacun de ces classifieurs de base. En particulier, nous ferons le lien entre cet algorithme et le principe de la minimisation du risque empirique énoncé par [Vapnik \(1999\)](#). Nous exposons en outre l'extension de cet algorithme au cas multi-classe.
- Nous exposons ensuite le cadre de l'apprentissage semi-supervisé dans le chapitre 7. Nous commençons ce chapitre par présenter les algorithmes EM et CEM développés dans le cadre de l'apprentissage non supervisé, en détaillant quelques cas particuliers menant à des modèles non supervisés bien connus

comme l'algorithme des k-moyennes. Nous présentons ensuite les hypothèses de base en apprentissage semi-supervisé en détaillant les trois approches génératives, discriminantes et graphiques développées suivant ce cadre.

- C'est dans le chapitre 8 que nous décrivons formellement le cadre de l'apprentissage de fonctions d'ordonnement (ou learning to rank en anglais) en focalisant sur deux formes particulières d'ordonnement appelées ordonnancement d'alternatives et d'ordonnement d'instances. Nous exposons ensuite quelques algorithmes développés suivant les approches classiques de l'apprentissage de fonctions d'ordonnement. Nous terminons ce chapitre par montrer la réduction de quelques problèmes d'ordonnement à la classification binaire de paires d'observations. Cette réduction ouvre la voie à l'apprentissage de classifieurs avec des exemples interdépendants que nous analysons avec le résultat de [Janson \(2004\)](#).
- En annexe A, nous donnons quelques rappels des outils de base en probabilité que nous employons dans cet ouvrage.
- En annexe B, nous donnons les codes programmes de quinze algorithmes présentés dans les différents chapitres, en détaillant les structures de données utilisées et en liant les différentes parties des programmes aux points correspondants abordés dans cet ouvrage.

Chapitre 1

Principes de base en apprentissage supervisé

Dans ce chapitre, nous exposons la théorie de l'apprentissage machine selon le cadre de [Vapnik \(1999\)](#) qui nous servira de base dans notre description des algorithmes d'apprentissage décrits dans les chapitres suivants. Plus particulièrement, nous présentons la notion de consistance qui garantit l'apprenabilité d'une fonction de prédiction. Les définitions et les hypothèses de base de cette théorie, ainsi que le principe de la minimisation du risque empirique, sont décrits dans la section 1.1. L'étude de la consistance de ce principe, présentée dans la section 1.2, nous mène au second principe de la minimisation du risque structural, qui stipule que l'apprentissage est un compromis entre une erreur empirique faible et une capacité de la classe de fonctions forte.

Un modèle d'apprentissage construit une fonction de prédiction à partir d'un ensemble fini d'exemples, appelé base d'entraînement ou base d'apprentissage (Fukunaga 1972 ; Duda *et al.* 2001 ; Schölkopf et Smola 2002 ; Boucheron *et al.* 2005). Suivant le cadre supervisé, chaque exemple est un couple constitué généralement du vecteur représentatif d'une observation et de sa réponse associée (aussi appelée sortie désirée). Le but de l'apprentissage est d'induire une fonction qui prédise les réponses associées à de nouvelles observations en commettant une erreur de prédiction la plus faible possible. Cette réponse est généralement une valeur réelle ou une étiquette de classe, comme nous allons le voir dans la suite. L'hypothèse sous-jacente ici est que les données sont stationnaires, c'est-à-dire que les exemples de la base d'entraînement, sur laquelle la fonction de prédiction est apprise, sont en quelque sorte représentatifs du problème général que l'on souhaite résoudre. Nous allons revenir sur cette hypothèse dans la section suivante.

En pratique, parmi une classe de fonctions existante, le modèle d'apprentissage choisit la fonction qui réalise la plus faible erreur moyenne de prédiction (ou erreur empirique) sur une base d'entraînement. La fonction d'erreur quantifie le désaccord entre la prédiction de sortie donnée par la fonction que l'on souhaite apprendre pour une observation de la base d'entraînement et sa réponse associée. Le but de cette recherche n'est pas que le modèle d'apprentissage induise une fonction donnant exactement les sorties désirées des observations de la base d'entraînement (ou faire du surapprentissage), mais de trouver, comme nous venons de l'évoquer, la fonction qui aura de bonnes performances de généralisation.

En logique, ce raisonnement ou procédé de recherche d'une règle générale à partir d'un ensemble d'observations fini est appelé induction (Genesereth et Nilsson 1987, chapitre 7, pp.161-176)¹. En apprentissage machine, le cadre inductif a été mis en place suivant le principe de la minimisation du risque empirique (MRE) (ou *Empirical Risk Minimisation* en anglais) et ses propriétés statistiques ont été étudiées dans la théorie développée par Vapnik (1999). Le résultat marquant de cette théorie est une borne supérieure de l'erreur de généralisation de la fonction apprise qui s'exprime en fonction de l'erreur empirique de cette dernière sur une base d'entraînement et de la complexité de la classe de fonctions utilisée. Cette complexité traduit la capacité de la classe de fonctions à résoudre le problème de prédiction et elle est d'autant plus grande qu'il y a de possibilités d'assigner des sorties désirées à des observations de la base d'entraînement. En d'autres termes, plus la capacité est grande, plus le risque empirique serait faible et moins on est garanti d'atteindre l'objectif principal de l'apprentissage, qui est d'avoir une faible erreur de généralisation. Cette borne exhibe ainsi le compromis qui existe entre l'erreur empirique et la capacité de la classe de fonctions, et montre une façon de

1. Le raisonnement contraire appelé déduction se base, quant à lui, sur des axiomes et produit des règles spécifiques (qui sont toujours vraies) comme des conséquences de la loi.

minimiser la borne sur l'erreur de généralisation (et d'avoir ainsi une meilleure estimation de cette erreur) en minimisant l'erreur empirique tout en contrôlant la capacité de l'ensemble de fonctions. Ce principe s'appelle la minimisation du risque structurel ; le principe ERM et lui sont à l'origine d'un grand nombre d'algorithmes d'apprentissage. De plus, ils peuvent expliquer le fonctionnement des algorithmes conçus avant l'établissement de la théorie de [Vapnik \(1999\)](#). La suite de ce chapitre est consacrée à la présentation plus formelle de ces différents concepts suivant le cadre de la classification bi-classe, qui a constitué le cadre initial du développement de cette théorie.

1.1 Principe de la minimisation du risque empirique

Dans cette section, nous allons présenter le principe de minimisation de risque empirique en fixant tout d'abord les notations qui seront utilisées par la suite.

1.1.1 Hypothèse et définitions

Nous supposons que les observations possèdent une représentation numérique dans un espace vectoriel de dimension fixe d , $\mathcal{X} \subseteq \mathbb{R}^d$. Les sorties désirées des observations sont supposées faire partie d'un ensemble de sortie $\mathcal{Y} \subset \mathbb{R}$. Jusqu'au début des années 2000, il y avait deux déclinaisons majeures des problèmes d'apprentissage supervisé ; la classification et la régression. En classification, l'ensemble de sortie \mathcal{Y} est discret et la fonction de prédiction $f : \mathcal{X} \rightarrow \mathcal{Y}$ est appelée un classifieur. Lorsque \mathcal{Y} est continu, f est une fonction de régression. Dans le chapitre 8, nous présenterons le cadre d'apprentissage de fonctions d'ordonnement qui s'est développé récemment dans les communautés de l'apprentissage machine et de la recherche d'information. Un couple $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ désigne ainsi un exemple étiqueté et $S = (\mathbf{x}_i, y_i)_{i=1}^m \in (\mathcal{X} \times \mathcal{Y})^m$ dénote un ensemble d'exemples d'entraînement. Dans le cas particulier de la classification binaire que l'on considère dans ce chapitre, nous notons l'espace de sortie par $\mathcal{Y} = \{-1, +1\}$ et un exemple $(\mathbf{x}, +1)$ (respectivement $(\mathbf{x}, -1)$) est appelé un exemple positif (respectivement négatif). Par exemple considérons le problème de classification de courriels, consistant à les étiqueter suivant deux classes : sollicité et non sollicité. On représentera les courriels par des vecteurs dans un espace vectoriel donné et on désignera une des classes (par exemple la classe des courriels sollicités) par l'étiquette de classe $+1$ et l'autre classe par l'étiquette de classe -1 .

L'hypothèse fondamentale de la théorie de l'apprentissage machine est que tous les exemples sont générés indépendamment et identiquement selon une distribution de probabilité fixe, mais inconnue, notée \mathcal{D} . L'hypothèse identiquement distribuée assure que les observations sont stationnaires, alors que l'hypothèse indépendamment distribuée stipule que chaque exemple individuel apporte un

maximum d'information pour résoudre le problème de prédiction. D'après cette hypothèse, les exemples (\mathbf{x}_i, y_i) de tout ensemble d'entraînement S et de test sont supposés être identiquement et indépendamment distribués (i.i.d.) selon \mathcal{D} . Autrement dit, chaque ensemble est un échantillon d'exemples i.i.d. selon \mathcal{D} .

Cette hypothèse caractérise ainsi la notion de représentativité d'un ensemble d'apprentissage et de test par rapport au problème de prédiction, c'est-à-dire que les exemples d'entraînement ainsi que les observations futures et leur sortie désirée sont supposés être issus d'une même source d'information.

Un autre concept de base en apprentissage est la notion de coût, aussi appelé risque ou erreur. Pour une fonction de prédiction f donnée, le désaccord entre la sortie désirée y d'un exemple \mathbf{x} et la prédiction $f(\mathbf{x})$ est mesurée grâce à une fonction de coût instantané définie par :

$$\mathbf{e} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$$

D'une manière générale, cette fonction est une distance sur l'ensemble de sortie \mathcal{Y} et elle mesure l'écart entre la réponse réelle et la réponse prédite par la fonction de prédiction pour une observation donnée. En régression, les fonctions de coût instantané usuelles sont les normes ℓ_1 et ℓ_2 de la différence entre les réponses réelle et prédite d'une observation donnée. En classification bi-classe, l'erreur instantanée communément envisagée est le coût 0/1, qui pour une observation (\mathbf{x}, y) et une fonction de prédiction f est définie par :

$$\mathbf{e}(f(\mathbf{x}), y) = \mathbb{1}_{f(\mathbf{x}) \neq y}$$

où $\mathbb{1}_\pi$ vaut 1 si le prédicat π est vrai et 0 sinon. En pratique, et dans le cas de la classification bi-classe, la fonction apprise $h : \mathcal{X} \rightarrow \mathbb{R}$ est une fonction à valeurs réelles et le classifieur associé $f : \mathcal{X} \rightarrow \{-1, +1\}$ est défini en prenant la fonction signe sur la sortie de h . Dans ce cas, l'erreur instantanée équivalente au coût 0/1, définie pour la fonction h est :

$$\begin{aligned} \mathbf{e}_0 : \mathbb{R} \times \mathcal{Y} &\rightarrow \mathbb{R}^+ \\ (h(\mathbf{x}), y) &\mapsto \mathbb{1}_{y \times h(\mathbf{x}) \leq 0} \end{aligned}$$

À partir d'un coût instantané et de la génération i.i.d. des exemples selon la distribution \mathcal{D} , on peut définir l'erreur de généralisation d'une fonction apprise $f \in \mathcal{F}$ comme :

$$\mathcal{L}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbf{e}(f(\mathbf{x}), y) = \int_{\mathcal{X} \times \mathcal{Y}} \mathbf{e}(f(\mathbf{x}), y) d\mathcal{D}(\mathbf{x}, y) \quad (1.1)$$

où $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} X(\mathbf{x}, y)$ est l'espérance de la variable aléatoire X lorsque (\mathbf{x}, y) suit la distribution de probabilité \mathcal{D} . Comme \mathcal{D} est inconnue, cette erreur de généralisation ne peut pas être estimée exactement, et pour mesurer la performance

d'une fonction f , on utilise souvent un ensemble d'exemples S de taille m sur lequel on calcule l'erreur empirique de f définie par :

$$\hat{\mathcal{L}}(f, S) = \frac{1}{m} \sum_{i=1}^m \mathbf{e}(f(\mathbf{x}_i), y_i) \quad (1.2)$$

Ainsi, pour résoudre un problème de classification pour lequel nous disposons d'un ensemble d'entraînement S , il est naturel de choisir une classe de fonctions \mathcal{F} et de chercher le classifieur f_S qui minimise l'erreur empirique sur S (puisque cette erreur est un estimateur non biaisé de l'erreur de généralisation de f_S que l'on ne peut pas mesurer).

1.1.2 Énoncé du principe

Cette méthode d'apprentissage, appelée le principe de minimisation du risque empirique (MRE), est à l'origine des tout premiers modèles d'apprentissage machine.

La question fondamentale qui se pose alors est : suivant le cadre MRE, *peut-on générer dans tous les cas une fonction de prédiction qui généralise bien à partir d'un ensemble d'observations fini* ? La réponse à cette question est bien évidemment non. Pour s'en convaincre, considérons le problème jouet de classification binaire suivant.

EXEMPLE Surapprentissage (Bousquet *et al.* 2003)

Supposons que la dimension d'entrée est $d = 1$. Prenons l'espace des observations \mathcal{X} ; l'intervalle $[a, b] \subset \mathbb{R}$ où a et b sont des réels tels que $a < b$ et l'espace des sorties est $\{-1, +1\}$. De plus, supposons que la distribution \mathcal{D} générant les couples d'exemples (\mathbf{x}, y) est une distribution uniforme sur $[a, b] \times \{-1\}$. Autrement dit, les exemples sont choisis de façon aléatoire sur l'intervalle $[a, b]$ et, pour chaque observation, la sortie désirée est -1 .

Considérons maintenant un algorithme d'apprentissage minimisant le risque empirique, en choisissant une fonction dans la classe des fonctions $\mathcal{F} = \{f : [a, b] \rightarrow \{-1, +1\}\}$ de la façon suivante; après avoir pris connaissance d'un ensemble d'apprentissage $S = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$ l'algorithme produit la fonction de prédiction f_S telle que :

$$f_S(\mathbf{x}) = \begin{cases} -1, & \text{si } \mathbf{x} \in \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \\ +1, & \text{sinon.} \end{cases}$$

Dans ce cas, le classifieur produit par l'algorithme d'apprentissage a un risque empirique égal à 0, et ceci pour n'importe quel ensemble d'apprentissage donné. Cependant, comme le classifieur fait une erreur sur tout l'ensemble infini $[a, b]$ sauf pour les exemples d'une base d'entraînement finie, de mesure nulle, son erreur de généralisation est toujours égale à 1.

1.2 Consistance du principe MRE

La question sous-jacente à la question précédente est : *dans quel cas le principe MRE est-il susceptible de générer une règle générale d'apprentissage ?* La réponse à cette question réside dans une notion statistique appelée consistance.

1.2.1 Définition

Ce concept indique les deux conditions qu'un algorithme d'apprentissage doit remplir, à savoir (a) l'algorithme doit renvoyer une fonction dont l'erreur empirique reflète son erreur de généralisation lorsque la taille de la base d'entraînement tend vers l'infini et, (b) dans le cas asymptotique, l'algorithme doit permettre de trouver une fonction qui minimise l'erreur de généralisation dans la classe de fonctions considérée. De façon formelle :

$$(a) \forall \epsilon > 0, \lim_{m \rightarrow \infty} \mathbb{P}(|\hat{\mathcal{L}}(f_S, S) - \mathcal{L}(f_S)| > \epsilon) = 0, \text{ noté, } \hat{\mathcal{L}}(f_S, S) \xrightarrow{\mathbb{P}} \mathcal{L}(f_S)$$

$$(b) \hat{\mathcal{L}}(f_S, S) \xrightarrow{\mathbb{P}} \inf_{g \in \mathcal{F}} \mathcal{L}(g)$$

Ces deux conditions impliquent ainsi la convergence en probabilité de l'erreur empirique $\hat{\mathcal{L}}(f_S, S)$ de la fonction de prédiction trouvée par l'algorithme d'apprentissage sur la base d'entraînement S , f_S , vers son erreur de généralisation $\mathcal{L}(f_S)$ et $\inf_{g \in \mathcal{F}} \mathcal{L}(g)$ (figure 1.1).

Une façon naturelle d'analyser la condition (a) de la consistance, exprimant le concept de la généralisation, est d'utiliser l'inégalité suivante :

$$|\mathcal{L}(f_S) - \hat{\mathcal{L}}(f_S, S)| \leq \sup_{g \in \mathcal{F}} |\mathcal{L}(g) - \hat{\mathcal{L}}(g, S)| \quad (1.3)$$

Nous voyons bien d'après cette inégalité qu'une condition suffisante pour généraliser est qu'asymptotiquement, l'erreur empirique de la fonction de prédiction, dont l'écart en valeur absolue entre cette erreur et son erreur de généralisation parmi toutes les autres fonctions d'une classe de fonctions \mathcal{F} donnée est la plus grande, tend vers l'erreur de généralisation de la fonction, soit :

$$\sup_{g \in \mathcal{F}} |\mathcal{L}(g) - \hat{\mathcal{L}}(g, S)| \xrightarrow{\mathbb{P}} 0 \quad (1.4)$$

Cette condition suffisante pour généraliser est une considération au pire cas et, d'après (équation 1.3), elle implique une convergence uniforme bilatérale pour toutes les fonctions de la classe \mathcal{F} . De plus, la condition (équation 1.4) ne dépend pas de l'algorithme considéré mais uniquement de la classe de fonctions \mathcal{F} . Ainsi, une condition nécessaire pour que le principe MRE soit consistant est que la classe de fonctions considérée soit restreinte (voir l'exemple sur le surapprentissage de la section précédente).