

Machine Learning

Massih-Reza Amini

*Université Grenoble Alpes, Grenoble INP, CNRS
Laboratoire d'Informatique de Grenoble*

Abstract

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'ordonnement, à un niveau accessible aux étudiants de master et aux élèves ingénieurs.

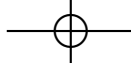
La première édition, connue sous le nom Apprentissage machine, fut traduite en chinois par les éditions iTuring. Dans cette deuxième édition, un nouveau chapitre est dédié au Deep Learning, sur les réseaux de neurones artificiels, et nous avons réorganisé les autres chapitres pour un exposé cohérent 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'ordonnement à large échelle.



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».

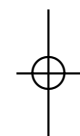
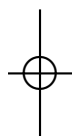
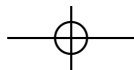
D'autres données ne sont pas si massives, comme les photos ou vidéos familiales, mais constituent toujours un défi algorithmique. Le grand changement récent est non seulement la taille, mais aussi le côté omniprésent de ces données, qui sont utilisées quotidiennement.

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



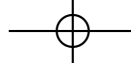
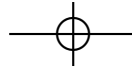


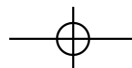
Table des matières

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	26
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



Machine Learning, de la théorie à la pratique

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	49
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



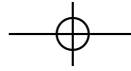
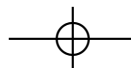
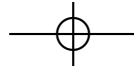


Table des matières

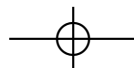
4.6.2	Méthode du moment	109
4.6.3	Traitement par mini-batches	109
4.6.4	Normalisation par batches	110
4.6.5	Technique de décrochage	111
CHAPITRE 5		
	Séparateurs à Vaste Marge	113
5.1	Notion de marge	114
5.1.1	Marge dure	114
5.1.2	Marge souple	119
5.2	Astuce du noyau	127
5.2.1	Définition de fonction 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	142
6.1.2	Échantillonnage par rejet	144
6.2	Étude théorique	144
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	149
6.3.1	Pseudo-erreur de classification	151
6.3.2	Échantillonnage par rejet suivant deux distributions	151
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





Machine Learning, de la théorie à la pratique

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



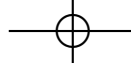
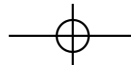
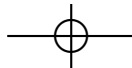


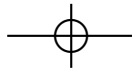
Table des matières

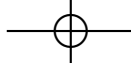
B.1 Structures de données	247
B.1.1 Base de données	247
B.1.2 Structure des hyper-paramètres	248
B.2 Structure pour une représentation creuse	249
B.3 Lancement des programmes	251
B.4 Codes	253
B.4.1 Algorithme BGFS (chapitre 3, section 3.2.2)	253
B.4.2 Recherche linéaire (chapitre 3, section 3.3)	256
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	303





Machine Learning, de la théorie à la pratique



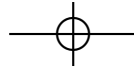


Avant-propos

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.



Machine Learning, de la théorie à la pratique

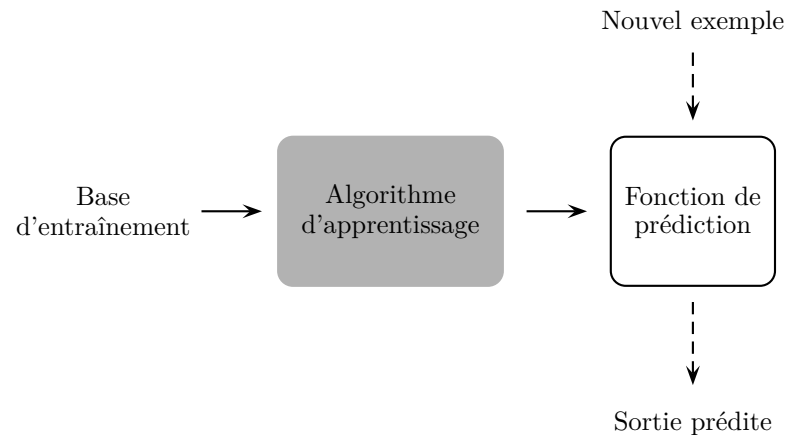
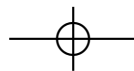


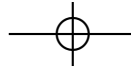
Figure 1 - Illustration des deux phases d'un problème d'apprentissage supervisé. Dans la phase d'apprentissage (schématisée par les traits pleins), une fonction minimisant l'erreur empirique sur une base d'entraînement est trouvée parmi une classe de fonctions prédéfinies. Dans la phase de test (schématisée 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 1). 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 (1). 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.

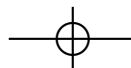
Depuis de nombreuses années, les algorithmes d'apprentissage développés suivant ces cadres ont été appliqués avec succès à une grande variété de problèmes, incluant la reconnaissance de la parole et de l'écriture manuscrite, la vision par ordinateur, la prédiction de la structure des protéines, les systèmes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

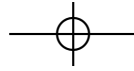
Organisation du livre

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 (1). 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>

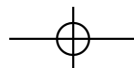


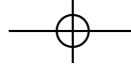


Machine Learning, de la théorie à la pratique

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. 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 (1). 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

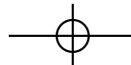


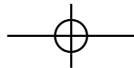
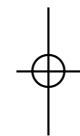
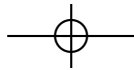


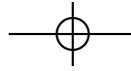
Avant-propos

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 (2).
- 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.

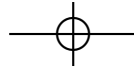






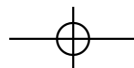
Bibliographie

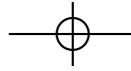
- [1] V. N. Vapnik, *The nature of statistical learning theory (second edition)*. Springer-Verlag, 1999.
- [2] S. Janson, “Large deviations for sums of partly dependent random variables,” *Random Structures and Algorithms*, vol. 24, no. 3, pp. 234–248, 2004.
- [3] M. Amini and N. Usunier, “A contextual query expansion approach by term clustering for robust text summarization,” in *Proceedings of the 7th Document Understanding Conference*, (Rochester - USA), 2007.
- [4] J.-F. Pessiot, Y.-M. Kim, M. Amini, and P. Gallinari, “Improving document clustering in a learned concept space,” *Information Processing & Management*, vol. 46, no. 2, pp. 180–192, 2010.
- [5] M. Amini, N. Usunier, and F. Laviolette, “A transductive bound for the voted classifier with an application to semi-supervised learning,” in *Advances in Neural Information Processing Systems 21*, pp. 65–72, 2009.
- [6] N. Usunier, M. Amini, and P. Gallinari, “A data-dependent generalisation error bound for the AUC,” in *ICML’05 workshop on ROC Analysis in Machine Learning*, 2005.
- [7] M. Amini, “Interactive learning for text summarization,” in *Proceedings of the PKDD/MLTIA Workshop on Machine Learning and Textual Information Access*, (Lyon - France), 2000.
- [8] Y.-M. Kim, M.-R. Amini, C. Goutte, and P. Gallinari, “Multi-view clustering of multilingual documents,” in *Proceedings of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR ’10, pp. 821–822, 2010.



Machine Learning, de la théorie à la pratique

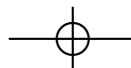
- [9] K. Fukunaga, *Introduction to Statistical Pattern Recognition*. New York, USA : Academic Press, 1972.
- [10] R. Duda, P. Hart, and D. Stork, *Pattern Classification*. Wiley, 2001.
- [11] B. Schölkopf and A. J. Smola, *Learning with kernels : support vector machines, regularization, optimization, and beyond*. MIT Press, 2002.
- [12] S. Boucheron, O. Bousquet, and G. Lugosi, “Theory of classification : a survey of some recent advances,” *ESAIM : Probability and Statistics*, pp. 323–375, 2005.
- [13] M. R. Genesereth and N. J. Nilsson, *Logical Foundations of Artificial Intelligence*. San Francisco, CA, USA : Morgan Kaufmann Publishers Inc., 1987.
- [14] O. Bousquet, S. Boucheron, and G. Lugosi, “Introduction to statistical learning theory,” in *Advanced Lectures on Machine Learning*, pp. 169–207, 2003.
- [15] J. Langford, “Tutorial on practical prediction theory for classification,” *Journal of Machine Learning Research*, vol. 6, pp. 273–306, Dec. 2005.
- [16] W. Hoeffding, “Probability inequalities for sums of bounded random variables,” *Journal of the American Statistical Association*, vol. 58, pp. 13–30, 1963.
- [17] P. L. Tchebychev, “Des valeurs moyennes,” *Journal de mathématiques pures et appliquées*, vol. 2, no. 12, pp. 177–184, 1867.
- [18] V. N. Vapnik and A. J. Chervonenkis, “On the uniform convergence of relative frequencies of events to their probabilities,” *Theory of Probability and its Applications*, vol. 16, pp. 264–280, 1971.
- [19] N. Sauer, “On the density of families of sets,” *Journal of Combinatorial Theory*, vol. 13, no. 1, pp. 145–147, 1972.
- [20] S. Shelah, “A combinatorial problem : Stability and order for models and theories in infinity languages,” *Pacific Journal of Mathematics*, vol. 41, pp. 247–261, 1972.
- [21] H. Brönnimann and M. T. Goodrich, “Almost optimal set covers in finite vc-dimension.,” *Discrete and Computational Geometry*, vol. 14, no. 4, pp. 463–479, 1995.
- [22] N. Cesa-Bianchi and D. Haussler, “A graph-theoretic generalization of the sauer-shelah lemma,” *Discrete Applied Mathematics*, vol. 86, pp. 27–35, 1998.

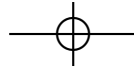




Bibliographie

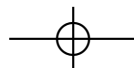
- [23] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of Machine Learning*. MIT Press, 2012.
- [24] V. N. Vapnik and A. J. Chervonenkis, “Theory of pattern recognition,” *Nauka*, 1974.
- [25] V. I. Koltchinskii and D. Panchenko, “Rademacher processes and bounding the risk of function learning,” in *High Dimensional Probability II* (E. Giné, D. Mason, and J. Wellner, eds.), pp. 443–459, 2000.
- [26] V. I. Koltchinskii, “Rademacher penalties and structural risk minimization,” *IEEE Transactions on Information Theory*, vol. 47, no. 5, pp. 1902–1914, 2001.
- [27] P. Massart, “Some applications of concentration inequalities to statistics,” *Annales de la faculté des sciences de Toulouse*, vol. 9, no. 2, pp. 245–303, 2000.
- [28] P. L. Bartlett and S. Mendelson, “Rademacher and gaussian complexities : risk bounds and structural results,” *Journal of Machine Learning Research*, vol. 3, pp. 463–482, 2003.
- [29] J. Taylor and N. Cristianini, *Kernel Methods for Pattern Analysis*. New York, USA : Cambridge Press University, 2004.
- [30] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth, “Learnability and the vapnik-chervonenkis dimension,” *Journal of the ACM*, vol. 36, pp. 929–965, 1989.
- [31] A. Ehrenfeucht, D. Haussler, M. Kearns, and L. Valiant, “A general lower bound on the number of examples needed for learning,” *Information and Computation*, vol. 82, pp. 247–261, 1989.
- [32] E. Giné, “Empirical processes and applications : an overview,” *Bernoulli*, vol. 2, no. 1, pp. 1–28, 1996.
- [33] C. McDiarmid, “On the method of bounded differences,” *Surveys in combinatorics*, vol. 141, pp. 148–188, 1989.
- [34] M. Ledoux and M. Talagrand, *Probability in Banach Spaces : Isoperimetry and Processes*. Springer Verlag, 1991.
- [35] A. Antos, B. Kégl, T. Linder, and G. Lugosi, “Data-dependent margin-based generalization bounds for classification,” *Journal of Machine Learning Research*, vol. 3, pp. 73–98, 2003.

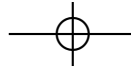




Machine Learning, de la théorie à la pratique

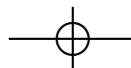
- [36] R. Kohavi, “A study of cross-validation and bootstrap for accuracy estimation and model selection,” in *Proceedings of the 14th International Joint Conference on Artificial Intelligence (IJCNN)*, pp. 1137–1143, 1995.
- [37] P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe, “Convexity, classification, and risk bounds,” *Journal of the American Statistical Association*, vol. 101, no. 473, pp. 138–156, 2006.
- [38] S. Boyd and L. Vandenberghe, *Convex Optimization*. New York, NY, USA : Cambridge University Press, 2004.
- [39] J. Nocedal and S. J. Wright, *Numerical Optimization*. Springer, 2006.
- [40] D. D. Lewis, Y. Yang, T. Rose, and F. Li, “RCV1 : A new benchmark collection for text categorization research,” *Journal of Machine Learning Research*, vol. 5, pp. 361–397, 2004.
- [41] F. Sebastiani, “Machine learning in automated text categorization,” *ACM Computing Surveys*, vol. 34, no. 1, pp. 1–47, 2004.
- [42] R. W. Hamming, “Error detecting and error correcting codes,” *Bell System Technical Journal*, vol. 29, no. 2, pp. 147–160, 1950.
- [43] T. G. Dietterich and G. Bakiri, “Solving multiclass learning problems via error-correcting output codes,” *Journal of Artificial Intelligence Research*, vol. 2, pp. 263–286, 1995.
- [44] E. L. Allwin, R. E. Schapire, and Y. Singer, “Reducing multiclass to binary : A unifying approach for margin classifiers,” *Journal of Machine Learning Research*, vol. 1, pp. 113 – 141, 2000.
- [45] Y. Guermeur, *SVM Multiclass, Théorie et Applications*. Habilitation à diriger des recherches, Université Nancy 1, 2007.
- [46] Y. Guermeur, “Sample complexity of classifiers taking values in \mathbb{R}^Q , application to multi-class SVMs,” *Communications in Statistics - Theory and Methods*, vol. 39, no. 3, pp. 543–557, 2010.
- [47] D. E. Rumelhart, G. E. Hinton, and R. Williams, “Learning internal representations by error propagation,” *Parallel Distributed Processing : Explorations in the Microstructure of Cognition*, vol. I, 1986.
- [48] L. Bottou, “Online algorithms and stochastic approximations,” in *Online Learning and Neural Networks* (D. Saad, ed.), Cambridge, UK : Cambridge University Press, 1998. revised, oct 2012.

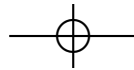




Bibliographie

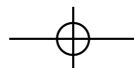
- [49] L. Bottou, “Large-scale machine learning with stochastic gradient descent,” in *Proceedings of the 19th International Conference on Computational Statistics (COMPSTAT’2010)* (Y. Lechevallier and G. Saporta, eds.), (Paris, France), pp. 177–187, Springer, August 2010.
- [50] F. Bach and E. Moulines, “Non-strongly-convex smooth stochastic approximation convergence rate $o(\frac{1}{n})$,” in *Advances in Neural Information Processing Systems (NIPS 26)*, pp. 773–781, 2013.
- [51] A. S. Nemirovski and D. B. Yudin, *Problem complexity and method efficiency in optimization*. Wiley-Interscience, 1983.
- [52] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, “Robust stochastic approximation approach to stochastic programming,” *Journal of the Society for Industrial and Applied Mathematics on Optimization (SIOPT)*, vol. 19, no. 4, pp. 1574–1609, 2009.
- [53] P. E. Gill and M. W. Leonard, “Reduced-hessian quasi-newton methods for unconstrained optimization,” *Journal of the Society for Industrial and Applied Mathematics on Optimization (SIOPT)*, vol. 12, no. 1, 2001.
- [54] P. Deuffhard, *Newton Methods for Nonlinear Problems : Affine Invariance and Adaptive Algorithms*. Springer Verlag, 2004.
- [55] J. F. Bonnans, J. C. Gilbert, C. Lemaréchal, and C. Sagastizàbal, *Numerical optimization, theoretical and numerical aspects*. Springer Verlag, 2006.
- [56] R. Fletcher, *Practical methods of optimization*. New York, USA : John Wiley & Sons, 1987.
- [57] W. C. Davidon, “Variable metric method for minimization,” *Journal of the Society for Industrial and Applied Mathematics on Optimization (SIOPT)*, vol. 1, no. 1, pp. 1–17, 1991.
- [58] E. Polak, *Computational methods in optimization*. Academic press, 1971.
- [59] P. Wolfe, “Convergence conditions for ascent methods,” *SIAM Review*, vol. 11, no. 2, pp. 226–235, 1966.
- [60] L. Armijo, “Minimization of functions having lipschitz continuous first partial derivatives,” *Pacific Journal of Mathematics*, vol. 16, no. 1, pp. 1–3, 1966.
- [61] G. Zoutendijk, “Some recent development in nonlinear programming,” in *5th Conference on Optimization Techniques, Part 1*, pp. 407–417, 1973.

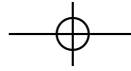




Machine Learning, de la théorie à la pratique

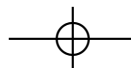
- [62] J. E. Dennis, Jr. and R. B. Schnabel, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations (Classics in Applied Mathematics, 16)*. Soc for Industrial & Applied Math, 1996.
- [63] K. A. Atkinson, *An introduction to numerical analysis*. John Wiley and Sons, 1988.
- [64] M. R. Hestenes and E. Stiefel, “Methods of conjugate gradients for solving linear systems,” *Journal of Research of the National Bureau of Standards*, vol. 49, pp. 409–436, 1952.
- [65] E. Polak and G. Ribiere, “Note sur la convergence de méthodes de directions conjuguées,” *ESAIM : Mathematical Modelling and Numerical Analysis - Modélisation Mathématique et Analyse Numérique*, vol. 3, no. R1, pp. 35–43, 1969.
- [66] R. Fletcher and C. M. Reeves, “Function minimization by conjugate gradients,” *The Computer Journal*, vol. 7, no. 2, pp. 149–154, 1964.
- [67] W. S. McCulloch and W. Pitts, “A logical calculus of the ideas immanent in nervous activity,” *Bulletin of Mathematical Biophysics*, vol. 5, pp. 115–133, 1943.
- [68] A. M. Turing, “Computing machinery and intelligence,” *Mind*, vol. 59, pp. 433–460, 1950.
- [69] D. Hebb, *The Organization of Behavior*. John Wiley, 1949.
- [70] M. Minsky, “A neural-analogue calculator based upon a probability model of reinforcement,” *Harvard University Psychological Laboratories, Cambridge, Massachusetts*, 1952.
- [71] F. Rosenblatt, “The perceptron : A probabilistic model for information storage and organization in the brain,” *Psychological Review*, vol. 65, pp. 386–408, 1958.
- [72] A. B. Novikoff, “On convergence proofs on perceptrons,” *Symposium on the Mathematical Theory of Automata*, vol. 12, pp. 615–622, 1962.
- [73] G. Widrow and M. Hoff, “Adaptive switching circuits,” *Institute of Radio Engineers, Western Electronic Show and Convention, Convention Record*, vol. 4, pp. 96–104, 1960.
- [74] N. J. Nilsson, *Learning machines; foundations of trainable pattern-classifying systems*. McGraw-Hill, 1965.
- [75] M. Minsky and S. Papert, *Perceptrons : An Introduction to Computational Geometry*. MIT Press, 1969.

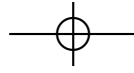




Bibliographie

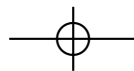
- [76] J. A. Anderson and E. Rosenfeld, *Neurocomputing : Foundations of Research*. MIT Press, 1988.
- [77] J. J. Hopfield, “Neurons with graded response have collective computational properties like those of two-state neurons,” in *Proceedings of the National Academy of Sciences USA*, pp. 3088–3092, 1984.
- [78] Y. LeCun, L. Bottou, and Y. Bengio, “Reading checks with multilayer graph transformer networks,” in *Proceedings of the 1997 IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP ’97)*, 1997.
- [79] S. Hocreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [80] V. N. Vapnik, *The nature of statistical learning theory*. 1995.
- [81] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems 25*, pp. 1097–1105, 2012.
- [82] Y. Freund and R. E. Schapire, “Large margin classification using the perceptron algorithm,” *Machine Learning Journal*, vol. 37, pp. 277–296, 1999.
- [83] J. Truett, J. Cornfield, and W. Kannel, “A multivariate analysis of the risk of coronary heart disease in framingham,” *Journal of Chronic Diseases*, vol. 20, no. 7, pp. 511–524, 1967.
- [84] J. A. Anderson, “Logistic discrimination,” *Handbook of Statistics*, vol. 2, pp. 169–191, 1982.
- [85] M. Kupperman, “Probabilities of hypotheses et information-statistics in sampling from exponential-class populations,” *Annals of Mathematical Statistics*, vol. 9, no. 2, pp. 571–575, 1958.
- [86] P. J. Werbos, *Beyond Regression : New Tools for Prediction and Analysis in the Behavioral Sciences*. PhD thesis, Harvard University, 1974.
- [87] D. Parker, “Learning logic,” Tech. Rep. TR-87, Cambridge, MA : Center for Computational Research in Economics and Management Science, MIT, 1985.
- [88] L. Bottou, *Une Approche théorique de l’Apprentissage Connexionniste : Applications à la Reconnaissance de la Parole*. PhD thesis, Université de Paris XI, Orsay, France, 1991.

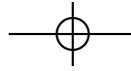




Machine Learning, de la théorie à la pratique

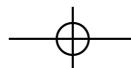
- [89] D. Hubel and T. Wiesel, “Receptive fields and functional architecture of monkey striate cortex,” *Journal of Physiology*, vol. 195, no. 1, pp. 215–243, 1968.
- [90] K. Fukushima, “Neocognitron : A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position,” *Biological Cybernetics*, vol. 36, no. 4, 1980.
- [91] A. Waibel, “Phoneme recognition using time-delay neural networks,” in *Meeting of the Institute of Electrical, Information and Communication Engineers*, 1987.
- [92] Y. LeCun, B. Boser, J. Denker, D. Henderson, R. Howard, W. Hubbard, and L. Jackel, “Backpropagation applied to handwritten zip code recognition,” *Neural Computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [93] A. J. Robinson and F. Fallside, “The utility driven dynamic error propagation network,” tech. rep., Engineering Department, Cambridge University, Cambridge, UK, 1987.
- [94] R. Williams and D. Zipser, “Gradient-based learning algorithms for recurrent networks and their computational complexity,” in *Backpropagation : theory, architectures, and applications* (B. Schölkopf, C. Burges, and A. Smola, eds.), ch. 4, p. 433–486, 1995.
- [95] M. D. Richard and R. P. Lippman, “Neural network classifiers estimate bayesian a posteriori probabilities,” *Neural Computation*, vol. 3, no. 4, pp. 461–483, 1991.
- [96] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*. Springer, 2001.
- [97] V. Nair and G. Hinton, “Rectified linear units improve restricted boltzmann machines,” in *Proceedings of the 27th International Conference on Machine Learning*, 2010.
- [98] B. Polyak, “Some methods of speeding up the convergence of iteration methods,” *USSR Computational Mathematics and Mathematical Physics*, vol. 4, no. 5, pp. 1–17, 1964.
- [99] Y. Nesterov, “A method of solving a convex programming problem with convergence rate $o(1/k^2)$,” *Soviet Mathematics Doklady*, vol. 27, pp. 372–376, 1983.
- [100] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, “On the importance of initialization and momentum in deep learning,” in *Proceedings of the 30th International Conference on Machine Learning*, pp. 1139–1147, 2013.

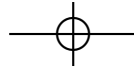




Bibliographie

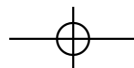
- [101] S. Ioffe and C. Szegedy, “Batch normalization : Accelerating deep network training by reducing internal covariate shift,” in *Proceedings of the 32nd International Conference on Machine Learning*, pp. 448–456, 2015.
- [102] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Improving neural networks by preventing co-adaptation of feature detectors,” *CoRR*, vol. abs/1207.0580, 2012.
- [103] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout : A simple way to prevent neural networks from overfitting,” *Journal of Machine Learning Research*, vol. 15, pp. 1929–1958, 2014.
- [104] B. Boser, I. Guyon, and V. Vapnik, “A training algorithm for optimal margin classifiers,” in *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, 1992.
- [105] T. Joachims, “Making large-scale SVM learning practical,” in *Advances in Kernel Methods - Support Vector Learning* (B. Schölkopf, C. Burges, and A. Smola, eds.), ch. 11, pp. 169–184, Cambridge, MA : MIT Press, 1999.
- [106] R. Fan, K. Chang, C. Hsieh, X. Wang, and C.-J. Lin, “LIBLINEAR : A library for large linear classification,” *Journal of Machine Learning Research*, vol. 9, pp. 1871–1874, 2008.
- [107] S. Shalev-Shwartz, Y. Singer, N. Srebro, and A. Cotter, “Pegasos : primal estimated sub-gradient solver for svm,” *Mathematical Programming*, vol. 127, no. 1, pp. 3–30, 2011.
- [108] J. Mercer, “Functions of positive and negative type and their connection with the theory of integral equations,” *Philosophical Transactions of the Royal Society*, vol. 209, pp. 415–446, 1909.
- [109] J. Weston and C. Watkins, “Support vector machines for multi-class pattern recognition,” in *European Symposium on Artificial Neural Networks (ESANN)*, pp. 219–224, 1999.
- [110] K. Crammer and Y. Singer, “On the algorithmic implementation of multi class kernel-based vector machines,” *Journal of Machine Learning Research*, vol. 2, pp. 265–292, 2001.
- [111] Y. Lee, Y. Lin, and G. Wahba, “Multicategory support vector machines : Theory and application to the classification of microarray data and satellite radiance data,” *Journal of the American Statistical Association*, vol. 99, no. 465, pp. 67–81, 2004.
- [112] L. Valiant, “The theory of the learnable,” *Communications of the ACM*, vol. 27, no. 11, pp. 1134–1142, 1984.

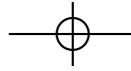




Machine Learning, de la théorie à la pratique

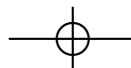
- [113] D. Kearns and L. Valiant, “Learning boolean formulae or finite automata is as hard as factoring,” Tech. Rep. TR-14-88, Harvard University Aiken Computation Laboratory, 1988.
- [114] R. E. Schapire, “The strength of weak learnability,” *Machine Learning*, vol. 5, no. 2, pp. 197–227, 1990.
- [115] R. E. Schapire, “Theoretical views of boosting and applications,” in *Proceedings of the 10th International Conference on Algorithmic Learning Theory*, pp. 13–25, 1999.
- [116] Y. Freund and R. E. Schapire, “A decision-theoretic generalization of on-line learning and an application to boosting,” *Journal of Computer and System Sciences*, vol. 55, no. 1, pp. 119–139, 1997.
- [117] R. E. Schapire, Y. Freund, P. Bartlett, and W. S. Lee, “Boosting the margin : a new explanation for the effectiveness of voting methods,” *The Annals of Statistics*, vol. 26, no. 5, pp. 1651–1680, 1998.
- [118] A. P. Dempster, N. M. Laird, and D. B. Rubin, “Maximum likelihood from incomplete data via the em algorithm,” *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 39, no. 1, pp. 1–38, 1977.
- [119] D. M. Titterton, A. F. M. Smith, and U. E. Smith, *Statistical Analysis of Finite Mixture Distributions*. Wiley, New York, 1985.
- [120] M. J. Symons, “Clustering criteria and multivariate normal mixtures,” *Biometrics*, vol. 37, no. 1, pp. 35–43, 1981.
- [121] G. Celeux and G. Govaert, “A classification em algorithm for clustering and two stochastic versions,” *Computational Statistics and Data Analysis*, vol. 14, no. 3, pp. 315–332, 1992.
- [122] T. Zhang and F. J. Oles, “A probability analysis on the value of unlabeled data for classification problems,” in *17th International Conference on Machine Learning*, 2000.
- [123] F. G. Cozman and I. Cohen, “Unlabeled data can degrade classification performance of generative classifiers,” in *Fifteenth International Florida Artificial Intelligence Society Conference*, pp. 327–331, 2002.
- [124] M. Seeger, “Learning with labeled and unlabeled data,” tech. rep., 2001.
- [125] S. Basu, A. Banerjee, and R. J. Mooney, “Semi-supervised clustering by seeding,” in *Proceedings of the Nineteenth International Conference on Machine Learning*, pp. 27 – 34, 2002.

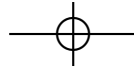




Bibliographie

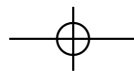
- [126] G. J. Machlachlan, *Discriminant Analysis and Statistical Pattern Recognition*. Wiley Interscience, 1992.
- [127] Y. Grandvalet and Y. Bengio, “Semi-supervised learning by entropy minimization,” in *Advances in Neural Information Processing Systems (NIPS 17)*, pp. 529–536, MIT Press, 2005.
- [128] O. Chapelle, B. Schölkopf, and A. Zien, eds., *Semi-Supervised Learning*. Cambridge, MA : MIT Press, 2006.
- [129] K. Nigam, A. K. McCallum, S. Thrun, and T. Mitchell, “Text classification from labeled and unlabeled documents using EM,” *Machine Learning Journal*, vol. 39, no. 2 - 3, pp. 103–134, 2000.
- [130] M.-R. Amini and É. Gaussier, *Recherche d’Information - applications, modèles et algorithmes*. Eyrolles, 2013.
- [131] I. Cohen, F. G. Cozman, N. Sebe, M. C. Cirelo, and T. S. Huang, “Semisupervised learning of classifiers : Theory, algorithms, and their application to human-computer interaction,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 26, no. 12, pp. 1553–1567, 2004.
- [132] S. C. Fralick, “Learning to recognize patterns without a teacher,” *IEEE Transactions on Information Theory*, vol. 13, no. 1, pp. 57–64, 1967.
- [133] E. A. Patrick, J. P. Costello, and F. C. Monds, “Decision-directed estimation of a two-class decision boundary,” *IEEE Transactions on Information Theory*, vol. 9, no. 3, pp. 197–205, 1970.
- [134] G. Tür, D. Z. Hakkani-Tür, and R. E. Schapire, “Combining active and semi-supervised learning for spoken language understanding,” *Speech Communication*, vol. 45, no. 2, pp. 171–186, 2005.
- [135] R. Urner, S. Shalev-Shwartz, and S. Ben-David, “Access to unlabeled data can speed up prediction time,” in *28th International Conference on Machine Learning, ICML 2011*, pp. 641–648, 2011.
- [136] P. Derbeko, E. El-Yaniv, and R. Meir, “Error bounds for transductive learning via compression and clustering,” in *Advances in Neural Information Processing Systems (NIPS 15)*, pp. 1085–1092, 2003.
- [137] T. Joachims, “Transductive inference for text classification using support vector machines,” in *Proceedings of the 16th International Conference on Machine Learning*, pp. 200–209, 1999.
- [138] T. Joachims, *Learning to Classify Text Using Support Vector Machines : Methods, Theory and Algorithms*. Norwell, MA, USA : Kluwer Academic Publishers, 2002.

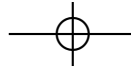




Machine Learning, de la théorie à la pratique

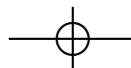
- [139] Z. Luo and P. Tseng, “On the convergence of the coordinate descent method for convex differentiable minimization,” *Journal of Optimization theory and applications*, vol. 72, no. 1, pp. 7–35, 1992.
- [140] M.-R. Amini, N. Usunier, and F. Laviolette, “A transductive bound for the voted classifier with an application to semi-supervised learning,” in *Advances in Neural Information Processing Systems (NIPS 21)*, pp. 65–72, 2009.
- [141] G. Dantzig, “Maximization of a linear function of variables subject to linear inequalities,” in *Activity Analysis of Production and Allocation* (T. Koopmans, ed.), p. 339–347, Wiley, New York, 1951.
- [142] F. Bach, R. Lanckriet, and M. Jordan, “Multiple kernel learning, conic duality, and the smo algorithm,” in *Proceedings of the Twenty-first International Conference on Machine Learning*, 2004.
- [143] V. Sindhwani, P. Niyogi, and M. Belkin., “A co-regularization approach to semi-supervised learning with multiple views,” in *ICML-05 Workshop on Learning with Multiple Views*, pp. 74–79, 2005.
- [144] A. Blum and T. Mitchell, “Combining labeled and unlabeled data with co-training,” in *Proceedings of the 11th Annual Conference on Learning Theory*, pp. 92–100, 1998.
- [145] B. Leskes, “The value of agreement, a new boosting algorithm,” in *Proceedings of Conference on Learning Theory (COLT)*, pp. 95–110, 2005.
- [146] X. Zhu and Z. Ghahramani, “Learning from labeled and unlabeled data with label propagation,” Tech. Rep. CMU-CALD-02-107, Carnegie Mellon University, 2002.
- [147] X. Zhu, Z. Ghahramani, and J. Lafferty, “Semi-supervised learning using gaussian fields and harmonic functions,” in *20th International Conference on Machine Learning*, pp. 912–919, 2003.
- [148] D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Schölkopf, “Learning with local and global consistency,” in *Advances in Neural Information Processing Systems (NIPS 16)*, pp. 321–328, MIT Press, 2004.
- [149] G. Latouche and V. Ramaswami, *Introduction to matrix analytic methods in stochastic modeling*. ASA-SIAM Series on Statistics and Applied Probability, Philadelphia, Pa. SIAM, Society for Industrial and Applied Mathematics Alexandria, Va. ASA, American Statistical Association, 1999.

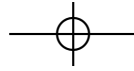




Bibliographie

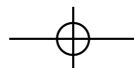
- [150] M. Szummer and T. Jaakkola, “Partially labeled classification with markov random walks,” in *Advances in Neural Information Processing Systems (NIPS 14)*, pp. 945–952, 2002.
- [151] E. Montroll, “Random walks in multidimensional spaces, especially on periodic lattices,” *Journal of the Society for Industrial and Applied Mathematics (SIAM)*, vol. 4, no. 4, pp. 241 – 260, 1956.
- [152] W. W. Cohen, R. E. Schapire, and Y. Singer, “Learning to order things,” in *Advances in Neural Information Processing Systems (NIPS 10)*, pp. 451–457, 1998.
- [153] C. Rudin, C. Cortes, M. Mohri, and R. E. Schapire, “Margin-based ranking meets boosting in the middle,” in *Conference On Learning Theory (COLT)*, 2005.
- [154] S. Agarwal, T. Graepel, R. Herbrich, S. Har-Peled, and D. Roth, “Generalization bounds for the area under the roc curve,” *Journal of Machine Learning Research*, vol. 6, pp. 393–425, 2005.
- [155] C. Cortes and M. Mohri, “AUC optimization vs. error rate minimization,” in *Advances in Neural Information Processing Systems (NIPS 16)*, pp. 313–320, 2004.
- [156] G. Salton, “A vector space model for automatic indexing,” *Communications of the ACM*, vol. 18, no. 11, pp. 613–620, 1975.
- [157] S. Hill, H. Zaragoza, R. Herbrich, and P. Rayner, “Average Precision and the Problem of Generalisation,” in *SIGIR Workshop on Mathematical and Formal Methods in Information Retrieval*, 2002.
- [158] S. E. Robertson and S. Walker, “Some simple effective approximations to the 2-poisson model for probabilistic weighted retrieval,” in *SIGIR’94, conference on Research and development in information retrieval*, pp. 232–241, 1994.
- [159] S. Clinchant and É. Gaussier, “Information-based models for *ad hoc* IR,” in *SIGIR’10, conference on Research and development in information retrieval*, 2010.
- [160] P. McCullagh, “Regression models for ordinal data,” *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 42, no. 2, pp. 109–142, 1980.
- [161] K. Crammer and Y. Singer, “Pranking with ranking,” in *Advances in Neural Information Processing Systems (NIPS 14)*, pp. 641–647, 2002.

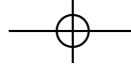




Machine Learning, de la théorie à la pratique

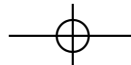
- [162] A. Shashua and A. Levin, “Ranking with large margin principle : Two approaches,” in *Advances in Neural Information Processing Systems (NIPS 15)*, pp. 961–968, 2003.
- [163] W. Chu and S. S. Keerthi, “New approaches to support vector ordinal regression,” in *22th International Conference on Machine Learning, ICML 2005*, pp. 145–152, 2005.
- [164] T. Qin, T.-Y. Liu, and H. Li, “A general approximation framework for direct optimization of information retrieval measures,” Tech. Rep. MSR-TR-2008- 164, Microsoft Research, November 2008.
- [165] M. Taylor, J. Guiver, S. Robertson, and T. Minka, “Softrank : Optimising non-smooth rank metrics,” in *WSDM 2008*, 2008.
- [166] Y. Yue, T. Finley, F. Radlinski, and T. Joachims, “A support vector method for optimizing average precision,” in *SIGIR ’07, conference on Research and Development in Information Retrieval*, pp. 271–278, 2007.
- [167] J. Xu and H. Li, “Adarank : A boosting algorithm for information retrieval,” in *SIGIR ’07, conference on Research and Development in Information Retrieval*, pp. 391 – 398, 2007.
- [168] J. Xu, T.-Y. Liu, M. Lu, H. Li, and W.-Y. Ma, “Directly optimizing evaluation measures in learning to rank,” in *SIGIR ’08, conference on Research and Development in Information Retrieval*, pp. 107–114, 2008.
- [169] C. Calauzènes, N. Usunier, and P. Gallinari, “On the (non-)existence of convex, calibrated surrogate losses for ranking,” in *Advances in Neural Information Processing Systems (NIPS 25)*, pp. 197–205, 2012.
- [170] Y. Freund, R. Iyer, R. E. Schapire, and Y. Singer, “An efficient boosting algorithm for combining preferences,” *Journal of Machine Learning Research*, vol. 4, pp. 933–969, 2003.
- [171] N. Usunier, *Apprentissage de fonctions d’ordonnement : une étude théorique de la réduction à la classification et deux applications à la Recherche d’Information*. PhD thesis, Université Pierre & Marie Curie, 2006.
- [172] S. K. Wong and Y. Y. Yao, “Linear structure in information retrieval,” in *SIGIR’88, conference on Research and Development in Information Retrieval*, pp. 219–232, 1988.
- [173] R. Herbrich, T. Graepel, P. Bollmann-Sdorra, and K. Obermayer, “Learning preference relations for information retrieval,” in *Proceedings of the AAAI Workshop Text Categorization and Machine Learning, Madison, USA.*, 1998.

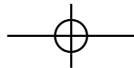
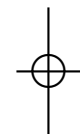
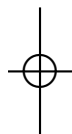
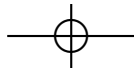




Index

- [174] A. Rakotomamonjy, “Optimizing area under roc curve with SVMs.,” in *1st International workshop on ROC Analysis in Artificial Intelligence*, pp. 71–80, 2004.
- [175] N. Usunier, M.-R. Amini, and P. Gallinari, “Generalization error bounds for classifiers trained with interdependent data,” in *Advances in Neural Information Processing Systems (NIPS 18)*, pp. 1369–1376, 2006.
- [176] H. Chernoff, “A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations,” *Annals of Mathematical Statistics*, vol. 23, no. 4, pp. 493–507, 1952.
- [177] P. Barbé and M. Ledoux, *Probabilité*. EDP Sciences, 2007.
- [178] W. Feller, *An Introduction to Probability Theory and Its Applications*. Wiley, 1968.
- [179] T. Bayes, “An essay towards solving a problem in the doctrine of chances,” *Philosophical Transactions of the Royal Society of London*, vol. 53, pp. 370–418, 1763.
- [180] P. S. Laplace, “Mémoire sur la probabilité des causes par les Événements,” *Académie Royale des sciences de Paris (Savants étrangers)*, vol. 6, pp. 621–656, 1771.
- [181] S. Boucheron, G. Lugosi, and P. Massart, *Concentration Inequalities : A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.





Index

- adaboost
 - cas bi-classe, 140, 209
 - algorithme, 141
 - code programme, 267
 - cas multi-classe, 150
 - algorithme, 153
 - code programme, 270
- adaline, 92
 - code programme, 261
- adaptive
 - boosting, *voir* adaboost
 - linear neuron, *voir* adaline
- affaiblissement du gradient, 107, 109
- aire sous la courbe ROC, 196
- alternatives, 200
- apprenant
 - faible, 140, 142, 146, 149–151
 - fort, 140
- apprentissage
 - multi-vue, 184
 - non supervisé, 155
 - semi-supervisé, 155
 - supervisé, 8
 - transductif, 176
- approximation quadratique, 56, 62, 75, 76
- area under curve, *voir* aire sous la courbe ROC
- astuce du noyau, 127
- AUC, *voir* aire sous la courbe ROC
- auto-apprenant, 174
 - code programme, 278
 - une passe, 176
- average pooling, *voir* regroupement moyen
- average precision, *voir* précision
- axiomes de Kolmogorov, 236
- back propagation, *voir*
 - rétro-propagation de l'erreur
- backpropagation through time, *voir* BPTT
- backtracking linesearch, *voir*
 - recherche linéaire stratégie retour en arrière
- bag of words, *voir* représentation sac de mots
- base
 - d'entraînement, 8, 11, 13, 14, 27, 34, 42, 43, 46–49, 51, 56, 60, 88, 90–92, 95, 97, 99–101, 114, 115, 118, 119, 121, 122, 125, 127, 130, 132, 137, 141, 145, 146, 148, 149, 151–153,

Machine Learning, de la théorie à la pratique

- 176, 177, 179, 180, 196, 207,
209, 210, 214–216, 220, 221,
223, 226, 229, 230, 232
- de test, 14–16
- de validation, 42
- boosting, 140
- borne
 - de généralisation, 17, 26, 50,
129, 130, 145, 148, 184, 220,
221, 226, 228
 - dépendante des données, 37
 - à base de marge, 39, 51,
129, 145, 148
 - de l’union, 17, 37, 237
 - de test, 14, 15, 219, 220
 - exponentielle, 144
 - uniforme, 14, 16
- BPTT, 107
- capacité d’une classe de fonctions,
8, 14, 22, 27, 28, 30
- changement covarié, 110
- classe de fonctions, 7, 8
- classification
 - binaire, 9–11
 - multi-classe
 - codes correcteurs, 49
 - mono-label, 45, 46, 50, 53
 - multi-label, 45, 46, 53
 - un contre tous, 46
 - un contre un, 47
- classification EM, 160
 - non supervisé
 - algorithme, 161
 - semi-supervisé
 - algorithme, 169
 - code programme, 275
- classification maximum likelihood,
voir maximum de vraisemblance
- classifiante
- classifieur, 9
 - de Gibbs, 180
 - de vote, 140, 141, 153, 179
 - faible, 140, 143
 - Naive Bayes, 170
- cluster assumption, *voir* hypothèse
de partition
- clustering, *voir* partitionnement
- co-training, 184
- coefficient multinomial, 171
- combinaison convexe, *voir*
enveloppe convexe
- complexité de Rademacher, 14, 30
 - empirique, 30, 33, 39, 52, 130
 - fractionnaire, 224–226
- condition
 - d’Armijo, 68, 69, 72
 - Karush–Kuhn–Tucker, 116,
120, 133
- condition de Mercer, 127
- conditions de Wolfe, 65, 69
- confiance, 130
- conjugaison mutuelle, 76, 77
- consistance, 12, 13, 21
- contrainte
 - non saturée, 115
 - saturée, 115
- contrainte de courbure, 68
- convergence
 - en probabilité, 12
- convergence uniforme
 - bilatérale, 12
 - unilatérale, 13
- convex hull, *voir* enveloppe convexe
convexe, 56
 - fortement, 124, 126
- convolutional neural network, *voir*
réseau convolutif
- coput
 - convexe, 56
- corps de Borel, *voir* tribu
- couche
 - de convolution, 103, 104, 106
 - de regroupement, 104–106
 - de sous-échantillonnage, 103
- courbe ROC, 196

- covariate shift, *voir* changement covarié
- coût
- 0/1, 10
 - logistique, 43, 97
 - convexe, 56, 61, 66
 - exponentiel, 44, 143
 - hinge, 43, 92, 123
 - quadratique, 92
 - à base de marge, 50
- crucial pair, *voir* paire cruciale
- DCG, *voir* gain cumulatif réduit
- decision boundary, *voir* frontière de décision
- decision stumps, 212
- densité, 240
- densité de probabilité, 156, 167
- dimension VC, 14, 21, 22, 26, 28, 30–32, 37, 53
- direction de descente, 61–63, 66, 67, 71–74, 77–79
- conjuguée, 75
 - gradient, 58
 - Newton, 61, 62
- discrimination logistique, *voir* régression logistique
- distance de Hamming, 45, 49
- distribution multinomiale, 171
- document frequency, 199
- dropout, *voir* décrochage
- dual de Wolfe, 117
- décision dirigée, *voir* auto-apprenant
- décrochage, 111
- déduction, 8
- développement de Taylor, 56, 59, 75, 245
- ECOC, *voir* classification multi-classe
- EM, *voir* Expectation Maximisation
- empirical risk minimization, *voir* minimisation du risque empirique
- ensemble de sortie, 9
- ensemble fondamental, 235
- enveloppe convexe, 39, 148
- erreur
- de classification, 44, 50, 130
 - de généralisation, 8, 10
 - empirique, 8, 11, 43
 - à base de marge, 50, 130, 132, 146–149
 - instantanée, 10, 38, 43
 - quadratique, 92
 - transductive, 177
- espace
- d'entrée, 9
 - de Hilbert, 127
 - de plongement, 127
 - de redescription, 119, 127, 130, 132, 146, 178
 - de représentation, 42
 - probabilisable, 236
 - probabilisé, 236, 237, 239
- espérance, 241
- estimateur non biaisé, 15, 35
- exemple étiqueté, 9
- Expectation Maximisation, 157
- algorithme, 159
- expérience aléatoire, 235
- feature space, *voir* espace de redescription, *voir* espace de redescription
- feed forward, 97
- filtre, 103
- fonction
- caractéristique, 86
 - d'activation, *voir* fonction de transfert
 - d'erreur, 8, 58
 - de croissance, 20, 25, 26, 30–33
 - de décision, 118
 - de projection, 127

Machine Learning, de la théorie à la pratique

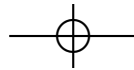
- de prédiction, 8–12, 14–16, 28
- de répartition, 240
- de seuillage, *voir* fonction de transfert
- de transfert, 99, 101, 108
- linéaire par morceaux, 108
- lipschitzienne, 39, 50, 130, 227
- objectif, 87, 92, 93, 116, 120, 122
- PReLU, 109
- ReLU, 108
- sigmoïde, 96, 108
- tangente hyperbolique, 108
- transfert, 85
- forme logistique, 95
- formule
 - BFGS, 63
 - de Bayes, 237
 - de Pascal, 25
 - de probabilités des causes, *voir* formule de Bayes
 - de Taylor-Lagrange, *voir* développement de Taylor
 - des probabilités composées, 237
 - des probabilités totales, 236, 238
 - du binôme, 25
 - Fletcher-Reeves, 79
 - Ribière-Polak, 78
- frontière de décision, 42
- gain cumulatif réduit, 197
- ghost sample, *voir* échantillon virtuel
- gradient
 - algorithme, 57, 124
 - conjugué, 73, 97
 - algorithme, 78
 - code programme, 258
 - descent, *voir* algorithme du gradient
 - lipschitzien, 71
 - stochastique, 61, 94
- growth function, *voir* fonction de croissance
- hard margin, *voir* marge dure
- hessienne, 56–60, 62, 63, 75, 76, 78, 116, 187
 - symétrique positive définie, 56, 63, 75
- hyperparamètre, 41
- hyperplan
 - marginal, 116–118, 121
 - séparateur, 86, 87, 89, 90, 94, 97, 114
- hypothèse
 - de continuité, 165, 187
 - de Naive Bayes, 170
 - de partition, 166, 167, 173
 - de variété, 166, 187
 - i.i.d., 10
- identiquement et indépendamment distribué, *voir* hypothèse i.i.d.
- idf, *voir* inverse document frequency
- induction, 8
- indépendants en probabilité, 239
- interpolation, 72, 73
- inégalité
 - de Bool, 237
 - de Cauchy-Schwarz, 38, 52, 131, 206, 223, 226
 - de Chernoff, 222, 243
 - de concentration, 34, 242
 - de Hoeffding, 15, 16, 21, 217, 219, 245
 - de Jensen, 31, 38, 40, 53, 126, 131, 158, 210, 223, 226, 244
 - de Markov, 243
 - de McDiarmid, 34–36
 - de Tchebychev, 20, 243
- jugement de pertinence, 194, 196, 197, 199–201, 203, 211

- k-moyennes
 - algorithmique, 163
 - code programme, 273
- kernel trick, *voir* astuce du noyau
- label spreading, *voir* propagation des étiquettes
- Lagrangien, 115, 116, 120, 133, 172
- learning rate, *voir* pas d'apprentissage
- learning to rank, *voir* ordonnancement
- leave- K -out cross-validation, *voir* validation croisée à K plis
- lemme
 - de Sauer, 22
 - de Massart, 31
 - de symétrisation, 19
 - de Talagrand, 39, 52, 131
- ligne de niveau, 57, 58
- likelihood, *voir* vraisemblance
- line search, *voir* recherche linéaire
- linéairement séparable, 89, 90, 94, 114, 117, 120, 123, 127, 132
- lissage de Laplace, 172
- listwise, *voir* ordonnancement approches
- logistic regression, *voir* régression
- loi d'une v.a.r. discrète, 241
 - Bernoulli, 241
- marche aléatoire, 190
- marge, 87, 90, 114, 118, 205, 207
 - ℓ_1 , 146
 - d'un exemple, 86
 - dure, 115
 - multi-classe, 50
 - souple, 119
- margin, *voir* marge
- matrice
 - de Gram, 128
 - définie positive, 64
 - jacobienne, 116
 - laplacienne, 187
 - par blocs, 187
 - stochastique, 189
- max pooling, *voir* regroupement maximal
- maximum de vraisemblance, 159, 161, 168
 - classifiante, 160, 161, 168, 171
- maximum likelihood, *voir* maximum de vraisemblance
- mean average precision, 195
- mesure de probabilité, 236
- mini-batch, 109, 110
- minimisation du risque
 - empirique, 8, 9, 11–13, 16, 46
 - structurel, 9, 14, 27, 28, 177
- mode
 - batch, 58, 99
 - en ligne, 60, 86, 90, 92
- modèle
 - informationnel, 202
 - OKAPI BM25, 202
 - vectorel, 198
- modèle vectoriel, 249
- mot de code, 49
- MRE, *voir* minimisation du risque empirique
- MRS, *voir* minimisation du risque structurel
- multi-layer perceptron, *voir* perceptron multi-couche
- mélange de densités, 156
- méthode de la sécante, 63
- méthodes d'ensemble, 140
- neocognitron, 103
- neurone formel, 84
- nombre chromatique, 219
 - pur, 219
- normalisation par lots, 110
- norme
 - ℓ_1 , 41
 - euclidienne, 41

Machine Learning, de la théorie à la pratique

- noyau
 - Gaussien, 128
 - polynomial, 128
 - positif, 127
 - sigmoïde, 128
 - symétrique, 127
- one vs all, *voir* classification multi-classe
- one vs one, *voir* classification multi-classe
- optimisation
 - d'ordre 1, 60
 - duale, 117, 122
 - sans contrainte, 55, 93
 - sous contraintes, 115, 121, 134
- ordonnement, 193
 - approches
 - par liste, 208
 - par paire, 208
 - par point, 203
 - instances
 - bipartite, 199, 211
- orthogonalité mutuelle, 77, 79
- outlier, *voir* point aberrant
- overfitting, *voir* surapprentissage
- paire cruciale, 195, 196, 199, 200, 208, 209, 211, 214, 216, 218, 219, 231
- pairwise, *voir* ordonnancement
 - approches
- partitionnement, 159
- pas d'apprentissage, 59, 88, 89, 94, 101, 102, 124, 248
 - adaptatif, 76
- pegasos, 124
 - algorithme, 125
- perceptron, 84, 204
 - algorithme, 88
 - code programme, 260
 - multi-couche, 97
 - algorithme de
 - rétro-propagation, 101
 - architecture, 98
 - code programme, 264
 - à marge, 91
- phase
 - propagation, 99, 111
 - rétropropagation, 99, 111
- PMC, *voir* perceptron multi-couche
- poids synaptiques, 84
- point aberrant, 121
- pointwise, *voir* ordonnancement
 - approches
- pooling, *voir* couche de regroupement
- PRank, 203
 - algorithme, 229
 - code programme, 282
- probabilité
 - P**-indépendants, *voir* indépendants en probabilité
 - σ -additivité, 236
 - σ -algèbre, *voir* tribu
 - a priori, 238
 - a posteriori, 238
 - conditionnelle, 237
 - écart-type, 241
 - échantillon
 - i.i.d., 10
 - virtuel, 19, 35
 - événement
 - aléatoire, 235
 - réalisable, 236, 237
 - événements incompatibles, 236–238
- programmation quadratique, 116
- propagation
 - de l'information, 99
 - des étiquettes, 186
- précision, 194
 - moyenne, 195
- pseudo-étiquette de classe, 176
- pulvérisation, 21

- qualification de contraintes, 115
- quasi-Newton, 61
 - algorithmme, 66
 - code programme, 253
 - condition, 63
- random walk, *voir* marche aléatoire
- rankboost, 209
 - algorithmme, 230
 - code programme, 284
- rappel, 194
- real time recurrent learning, *voir*
 - RTRL
- recall, *voir* rappel
- Receiver Operating Characteristics, *voir* courbe ROC
- recherche linéaire, 66, 67, 72–75, 179
 - algorithmme, 74
 - code programme, 256
 - stratégie de retour en arrière, 72
- recouvrement, 218
 - fractionnaire, 218, 225
 - propre exact, 218, 219, 221, 222
- regroupement
 - maximal, 105, 106
 - moyen, 105
- représentation creuse, 249
- représentation sac de mots, 198
- risque, 10
 - transductif, 180
- RKHS, 129
- routage d'information, 193
- RTRL, 107
- règle
 - de Bayes, 160
 - de décision bayésienne, 159, 160
 - de la chaîne, 102
- régression, 9
 - linéaire, 92
 - logistique, 95
 - code programme, 262
 - multi-classe, 108
 - ordinaire, 203
- régularisation, 41
- réseaux
 - convolutifs, 83, 103, 106
 - profonds, 98
 - récurrents, 97, 106, 107
 - sans boucle, 97
 - à retard temporel, 103
- rétro-propagation de l'erreur, 99, 104
 - algorithmme, 100
- self-training, *voir* auto-apprenant
- slack variables, *voir* variables
- smoothness assumption, *voir*
 - hypothèse de continuité
- soft margin, *voir* marge souple
- somme des carrés des résidus, 162
- sortie désirée, 8
- steepest descent, *voir* algorithmme
 - du gradient
- structural risk minimization, *voir*
 - minimisation du risque structurel
- structure de données, 248
- sum of squared residuals, *voir*
 - somme des carrés des résidus
- support vector machine, *voir*
 - séparateurs à vaste marge
- support vectors, *voir* vecteurs de support
- surapprentissage, 8, 11, 111
- symbole de Kronecker, 133
- symétrisation, 19, 35
- système complet d'événements, 236
- séparateurs à vaste marge, 113
 - algorithmme (marge dure), 119
 - algorithmme (marge souple), 122
 - multi-classe, 53, 132
 - algorithmme, 137
 - transductifs, 178
- taux de décroissance, 68



Machine Learning, de la théorie à la pratique

- technique d'échantillonnage par
 - rejet, 144, 152
- term frequency, 199
- tf-idf, 199
- théorème
 - de Hoeffding, *voir* inégalité de Hoeffding
 - de Janson, 217
 - de Novikoff, 89, 132
 - de Schwarz, 56
 - de Vapnik et Chervonenkis, 20
 - de Zoutendijk, 71
 - des différences bornées, *voir* inégalité de McDiarmid
 - des gendarmes, 126
- time delay neural network, *voir* réseaux à retard temporel
- training set, *voir* base d'entraînement
- transfert
 - Heaviside, 85
 - linéaire, 86
- tribu, 236
 - des boréliens, 239, 240
- univers des possibles, *voir* ensemble fondamental
- v.a.r., *voir* variables aléatoires réelles
- valeur propre, 56, 57, 60
- valeurs d'utilité, 194
- validation croisée à K plis, 42
- vanishing gradient, *voir* affaiblissement du gradient
- variables
 - aléatoires interdépendantes, 217, 221
 - aléatoires réelles, 239
 - mesurable, 240
 - d'écart, 119, 120, 179
 - de Lagrange, 115, 120, 122, 171
 - de Rademacher, 30, 131
 - discrètes, 240
 - primales, 116
- variables aléatoires
 - à densité, 240
- variance, 241
- vecteur
 - gradient, 56, 58, 60, 62, 74, 76, 77
 - indicateur
 - de classe, 45, 100, 165, 168, 172
 - de groupe, 160, 168
 - orthonormal, 135
 - propre, 56–59
- vecteur indicateur
 - de classe, 134
- vecteur support, 123
- vecteurs de support, 114, 116–118, 120, 121
- vector space model, 198
- vitesse de convergence, 60
 - sous-linéaire, 60
- vocabulaire, 250
- vraisemblance, 156
 - classifiante, 95, 160, 168, 171
- weak learner, *voir* apprenant faible
- worst case study, *voir* étude de pire cas
- étude de pire cas, 13

