



N°d'ordre NNT : ?

THÈSE de DOCTORAT DE L'UNIVERSITÉ DE LILLE

École Doctorale MADIS N° 631

**Mathématiques-Sciences du numérique et de leurs
interactions**

Spécialité / discipline de doctorat : Mathématique

Soutenue par :

Maxime Haddouche

A PAC-Bayes Approach of Generalisation In Machine Learning: From Heavy-Tailed Martingales To Optimisation

ACKNOWLEDGEMENTS

TODO

The further backward you look, the further forward you can see.

WINSTON CHURCHILL

Science is a bit like the joke about the drunk who is looking under a lamppost for a key that he has lost on the other side of the street, because that's where the light is. It has no other choice.

NOAM CHOMSKY

CONTENTS

Contents	5
List of Figures	6
List of Theorems	8
List of Notations	9
List of Publications	11
International Conference	11
International Workshop	11
National Conference	11
Research Report	12
 I Background	 13
 II PAC-Bayesian Majority Vote: Theory and Self-bounding Algorithms	 15
 III Beyond PAC-Bayesian Bounds: From Disintegration to Novel Bounds	 17
 IV Conclusion and Perspectives	 19
 V Appendix	 21
 A Some Mathematical Tools	 23
A.1 JENSEN's Inequality	23
A.2 MARKOV's Inequality	23
A.3 2nd Order MARKOV's Inequality	24
A.4 CHEBYSHEV-CANTELLI Inequality	25
A.5 HÖLDER's Inequality	25

LIST OF FIGURES

LIST OF ALGORITHMS

LIST OF THEOREMS

Appendix

A.1.1 Theorem (JENSEN's Inequality)	23
A.2.1 Theorem (MARKOV's Inequality)	23
A.3.1 Theorem (2nd Order MARKOV's Inequality)	24
A.4.1 Theorem (CHEBYSHEV-CANTELLI Inequality)	25
A.5.1 Lemma (YOUNG's Inequality)	25
A.5.1 Theorem (HÖLDER's Inequality)	26

LIST OF NOTATIONS

General

a	A scalar (integer or real)
\mathbf{a}	A vector
\mathbf{A}	A matrix
\mathbb{A}, \mathfrak{A}	A set
\mathbb{R}	The set of real numbers
\mathbb{R}_*	The set of real numbers excluding 0
\mathbb{R}_*^+	The set of positive real numbers excluding 0
\mathbb{N}	The set of natural numbers
\mathbb{N}_*	The set of natural numbers excluding 0
$\text{card}(\cdot)$	The cardinal of a set
a_i	i -th element of the vector \mathbf{a}

Statistical Learning Theory

\mathbb{X}	Set of d -dimensional inputs ($\subseteq \mathbb{R}^d$)
\mathbb{Y}	Set of labels
\mathbf{x}	A real-valued input $\mathbf{x} \in \mathbb{X}$
y	A label $y \in \mathbb{Y}$ associated to the input \mathbf{x}
\mathcal{D}	Unknown data distribution on $\mathbb{X} \times \mathbb{Y}$
\mathcal{D}^m	Unknown data distribution on the m -samples, <i>i.e.</i> , on $(\mathbb{X} \times \mathbb{Y})^m$
\mathbb{S}	Learning sample $\mathbb{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ drawn from \mathcal{D}^m
\mathcal{S}	Uniform distribution on \mathbb{S}
\mathbb{T}	Test set drawn from \mathcal{D}^m

\mathcal{T}	Uniform distribution on \mathbb{T}
\mathbb{H}	The set of hypotheses
h	A hypothesis $h \in \mathbb{H}$
$\ell(\cdot, \cdot)$	Loss function
$R_{\mathcal{D}'}^\ell(h)$	Risk of the hypothesis $h \in \mathbb{H}$ w.r.t. the loss function $\ell(\cdot)$ on \mathcal{D}'
$R_{\mathcal{D}'}(h)$	Risk of the hypothesis $h \in \mathbb{H}$ w.r.t. the 0-1 loss on \mathcal{D}'

Probability Theory

$\mathbb{E}_{X \sim \mathcal{X}}[\cdot]$	The expectation w.r.t. the random variable $X \sim \mathcal{X}$
$\mathbb{P}_{X \sim \mathcal{X}}[\cdot]$	The probability w.r.t. the random variable $X \sim \mathcal{X}$
$\mathbb{I}[a]$	Indicator function; returns 1 if a is true and 0 otherwise
$\mathbb{M}(\mathbb{H})$	Set of Probability densities w.r.t. the reference measure on \mathbb{H}
ρ	Posterior distribution $\rho \in \mathbb{M}(\mathbb{H})$ on \mathbb{H}
π	Prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$ on \mathbb{H}
$\text{KL}(\rho \parallel \pi)$	Kullback-Leibler (KL) divergence between ρ and π
$D_\alpha(\rho \parallel \pi)$	Rényi Divergence between ρ and π
$\text{Uni}(\mathbb{A})$	Uniform distribution on \mathbb{A}
$\text{Dir}(\boldsymbol{\alpha})$	Dirichlet distribution of parameters $\boldsymbol{\alpha} \in \mathbb{R}_*^+$

Majority Vote

$\text{MV}_\rho(\cdot)$	The majority vote classifier
$m_\rho(\mathbf{x}, y)$	Majority vote's margin for the example $(\mathbf{x}, y) \in \mathbb{X} \times \mathbb{Y}$
$\widehat{m}_\rho(\mathbf{x}, y)$	$\frac{1}{2}$ -Margin for the example $(\mathbf{x}, y) \in \mathbb{X} \times \mathbb{Y}$
$\text{sign}(a)$	Sign function; returns $+1$ if $a \geq 0$ and -1 otherwise
$r_{\mathcal{D}'}(\rho)$	Gibbs risk on the distribution \mathcal{D}' associated to the majority vote $\text{MV}_\rho(\cdot)$
$e_{\mathcal{D}'}(\rho)$	Joint Error on the distribution \mathcal{D}' associated to the majority vote $\text{MV}_\rho(\cdot)$
$d_{\mathcal{D}'}(\rho)$	Disagreement on the distribution \mathcal{D}' associated to the majority vote $\text{MV}_\rho(\cdot)$

LIST OF PUBLICATIONS

International Conference

PAUL VIALARD, PASCAL GERMAIN, AMAURY HABRARD, and EMILIE MORVANT. Self-bounding Majority Vote Learning Algorithms by the Direct Minimization of a Tight PAC-Bayesian C-Bound. *Machine Learning and Knowledge Discovery in Databases (ECML PKDD)*. (2021a).

PAUL VIALARD, GUILLAUME VIDOT, AMAURY HABRARD, and EMILIE MORVANT. A PAC-Bayes Analysis of Adversarial Robustness. *Advances in Neural Information Processing Systems (NeurIPS)*. (2021d).

VALENTINA ZANTEDESCHI, PAUL VIALARD, EMILIE MORVANT, RÉMI EMONET, AMAURY HABRARD, PASCAL GERMAIN, and BENJAMIN GUEDJ. Learning Stochastic Majority Votes by Minimizing a PAC-Bayes Generalization Bound. *Advances in Neural Information Processing Systems (NeurIPS)*. (2021).

International Workshop

PAUL VIALARD, PASCAL GERMAIN, AMAURY HABRARD, and EMILIE MORVANT. Interpreting Neural Networks as Majority Votes through the PAC-Bayesian Theory. *NeurIPS 2019 Workshop on Machine Learning with Guarantees*. (2019).

National Conference

PAUL VIALARD, RÉMI EMONET, PASCAL GERMAIN, AMAURY HABRARD, EMILIE MORVANT, and VALENTINA ZANTEDESCHI. Intérêt des bornes désintégrées pour la généralisation avec des mesures de complexité. *Conférence sur l'Apprentissage automatique (CAp)*. (2022a).

VALENTINA ZANTEDESCHI, PAUL VIALARD, EMILIE MORVANT, RÉMI EMONET, AMAURY HABRARD, PASCAL GERMAIN, and BENJAMIN GUEDJ. Learning Stochastic Majority Votes by Minimizing a PAC-Bayes Generalization Bound. *Conférence sur l'Apprentissage automatique (CAp)*. (2022).

PAUL VIALARD, PASCAL GERMAIN, and EMILIE MORVANT. Apprentissage de Vote de Majorité par Minimisation d'une C-Borne PAC-Bayésienne. *Conférence sur l'Apprentissage automatique (CAp)*. (2021b).

PAUL VIALARD, PASCAL GERMAIN, and EMILIE MORVANT. Dérandomisation des Bornes PAC-Bayésiennes. *Conférence sur l'Apprentissage automatique (CAp)*. (2021c).

GUILLAUME VIDOT, PAUL VIALARD, and EMILIE MORVANT. Une Analyse PAC-Bayésienne de la Robustesse Adversariale. *Conférence sur l'Apprentissage automatique (CAp)*. (2021).

PAUL VIALARD, RÉMI EMONET, AMAURY HABRARD, EMILIE MORVANT, and PASCAL GERMAIN. Théorie PAC-Bayésienne pour l'apprentissage en deux étapes de réseaux de neurones. *Conférence sur l'Apprentissage automatique (CAp)*. (2020).

Research Report

PAUL VIALARD, RÉMI EMONET, AMAURY HABRARD, EMILIE MORVANT, and VALENTINA ZANTEDESCHI. Generalization Bounds with Arbitrary Complexity Measures. *Submitted to ICLR 2023*. (2022b).

PAUL VIALARD, PASCAL GERMAIN, AMAURY HABRARD, and EMILIE MORVANT. A General Framework for the Disintegration of PAC-Bayesian Bounds. *Submitted to Machine Learning Journal*. (2022c).

PART I

Background

PART II

PAC-Bayesian Majority Vote: Theory and Self-bounding Algorithms

PART III

Beyond PAC-Bayesian Bounds: From Disintegration to Novel Bounds

PART IV

Conclusion and Perspectives

PART V

Appendix

SOME MATHEMATICAL TOOLS

A.1 Jensen's Inequality

Theorem A.1.1 (JENSEN's Inequality). Let $X \in \mathbb{X}$ a random variable following a probability distribution \mathcal{X} with $f : \mathbb{X} \rightarrow \mathbb{R}$ a measurable convex function, we have

$$f\left(\mathbb{E}_{X \sim \mathcal{X}}[X]\right) \leq \mathbb{E}_{X \sim \mathcal{X}}[f(X)].$$

Proof. Since $f()$ is a convex function, the following inequality holds, i.e., we have

$$\forall X' \in \mathbb{X}, \quad a\left(X' - \mathbb{E}_{X \sim \mathcal{X}}[X]\right) \leq f(X') - f\left(\mathbb{E}_{X \sim \mathcal{X}}[X]\right),$$

where a is the tangent's slope. By taking the expectation to both sides of the inequality, we have

$$\underbrace{a\left(\mathbb{E}_{X \sim \mathcal{X}}[X] - \mathbb{E}_{X \sim \mathcal{X}}[X]\right)}_{=0} \leq \mathbb{E}_{X \sim \mathcal{X}}[f(X)] - f\left(\mathbb{E}_{X \sim \mathcal{X}}[X]\right).$$

Hence, by rearranging the terms, we prove the claimed result. ■

A.2 Markov's Inequality

Theorem A.2.1 (MARKOV's Inequality). Let $X \in \mathbb{X}$ a non-negative random variable following a probability distribution \mathcal{X} and $\tau > 0$, we have

$$\mathbb{P}_{X \sim \mathcal{X}}[X \geq \tau] \leq \frac{\mathbb{E}_{X \sim \mathcal{X}}[X]}{\tau}.$$

Proof. First of all, remark that we have the following inequality for any $X \in \mathbb{X}$

$$\tau \mathbb{I}[X \geq \tau] \leq X \mathbb{I}[X \geq \tau] \leq X. \quad (\text{A.1})$$

Indeed, on the one hand, if $X < \tau$, $\mathbb{I}[X \geq \tau] = 0$, the inequality holds trivially. On the other hand, if $X \geq \tau$, $\mathbb{I}[X \geq \tau] = 1$ and the inequality becomes $\tau \leq X$, which is true. By taking the expectation of Equation (A.1), we have

$$\mathbb{E}_{X \sim \mathcal{X}} [\tau \mathbb{I}[X \geq \tau]] \leq \mathbb{E}_{X \sim \mathcal{X}} [X].$$

From the fact that the expectation of a constant is the constant and by definition of the probability, we have

$$\tau \mathbb{P}_{X \sim \mathcal{X}} [X \geq \tau] \leq \mathbb{E}_{X \sim \mathcal{X}} [X] \iff \mathbb{P}_{X \sim \mathcal{X}} [X \geq \tau] \leq \frac{\mathbb{E}_{X \sim \mathcal{X}} [X]}{\tau},$$

which is the desired result. ■

A.3 2nd Order Markov's Inequality

Theorem A.3.1 (2nd Order MARKOV's Inequality). Let X a non-negative random variable following a probability distribution \mathcal{X} and $\tau > 0$, we have

$$\mathbb{P}_{X \sim \mathcal{X}} [X \geq \tau] \leq \frac{\mathbb{E}_{X \sim \mathcal{X}} [X^2]}{\tau^2}.$$

Proof. We apply MARKOV's inequality (Theorem A.2.1) to have

$$\mathbb{P}_{X \sim \mathcal{X}} [X^2 \geq \tau^2] \leq \frac{\mathbb{E}_{X \sim \mathcal{X}} [X^2]}{\tau^2}.$$

Moreover, since $\mathbb{I}[X \geq \tau] = \mathbb{I}[X^2 \geq \tau^2]$, we have

$$\mathbb{P}_{X \sim \mathcal{X}} [X \geq \tau] = \mathbb{P}_{X \sim \mathcal{X}} [X^2 \geq \tau^2],$$

which proves the desired result. ■

A.4 Chebyshev-Cantelli Inequality

Theorem A.4.1 (CHEBYSHEV-CANTELLI Inequality). Let X a random variable following a probability distribution \mathcal{X} and $\tau > 0$, we have

$$\mathbb{P}_{X \sim \mathcal{X}} \left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' \geq \tau \right] \leq \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\mathbb{V}_{X' \sim \mathcal{X}} X' + \tau^2}.$$

Proof. First of all, remark that we have

$$\begin{aligned} \mathbb{P}_{X \sim \mathcal{X}} \left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' \geq \tau \right] &= \mathbb{P}_{X \sim \mathcal{X}} \left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \geq \tau + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \right] \\ &\leq \mathbb{P}_{X \sim \mathcal{X}} \left[\left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \right]^2 \geq \left[\tau + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \right]^2 \right], \end{aligned} \quad (\text{A.2})$$

where $\mathbb{V}_{X \sim \mathcal{X}} X$ is the variance of the random variable $X \sim \mathcal{X}$. From Equation (A.2) and MARKOV's Inequality (Theorem A.2.1), we can deduce that

$$\begin{aligned} \mathbb{P}_{X \sim \mathcal{X}} \left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' \geq \tau \right] &\leq \frac{\mathbb{E}_{X \sim \mathcal{X}} \left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \right]^2}{\left[\tau + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \right]^2} \\ &= \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\mathbb{V}_{X' \sim \mathcal{X}} X' + \tau^2}. \end{aligned}$$

■

A.5 Hölder's Inequality

In order to prove HÖLDER's inequality, we first prove the following lemma.

Lemma A.5.1 (YOUNG's Inequality). For any $\alpha > 1$ and $\beta > 1$ such that $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, for any $a \geq 0$ and $b \geq 0$, we have

$$ab \leq \frac{a^\alpha}{\alpha} + \frac{b^\beta}{\beta}.$$

Proof. We first develop $\ln [ab]$ and we apply JENSEN's inequality (Theorem A.1.1) since the logarithm is concave and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$. Indeed, we have

$$\ln [ab] = \ln a + \ln b = \frac{\alpha}{\alpha} \ln a + \frac{\beta}{\beta} \ln b = \frac{\ln a^\alpha}{\alpha} + \frac{\ln b^\beta}{\beta} \leq \ln \left[\frac{a^\alpha}{\alpha} + \frac{b^\beta}{\beta} \right].$$

Then, we take the exponential to both sides of the inequality and we are done. ■

We are now ready to prove HÖLDER's inequality.

Theorem A.5.1 (HÖLDER's Inequality). For any measurable function $f()$ and $g()$, for any $\alpha > 1$ and $\beta > 1$ such that $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, we have

$$\mathbb{E}_{X \sim \mathcal{X}} |f(X)g(X)| \leq \left[\mathbb{E}_{X \sim \mathcal{X}} |f(X)|^\alpha \right]^{\frac{1}{\alpha}} \left[\mathbb{E}_{X \sim \mathcal{X}} |g(X)|^\beta \right]^{\frac{1}{\beta}}.$$

Proof. For convenience of notation, let $\|f\|_\alpha = [\mathbb{E}_{X \sim \mathcal{X}} |f(X)|^\alpha]^{\frac{1}{\alpha}}$ and $\|g\|_\beta = [\mathbb{E}_{X \sim \mathcal{X}} |g(X)|^\beta]^{\frac{1}{\beta}}$. If $\|f\|_\alpha = 0$ or $\|g\|_\beta = 0$, then $\mathbb{E}_{X \sim \mathcal{X}} |f(X)g(X)| = 0$, hence, the inequality holds in this case. Then for $\|f\|_\alpha > 0$ and $\|g\|_\beta > 0$, we upper-bound the term $\frac{|f(X)g(X)|}{\|f\|_\alpha \|g\|_\beta}$ with YOUNG's inequality (Lemma A.5.1), i.e., we have

$$\frac{|f(X)g(X)|}{\|f\|_\alpha \|g\|_\beta} \leq \frac{|f(X)|^\alpha}{\alpha \|f\|_\alpha^\alpha} + \frac{|f(X)|^\beta}{\beta \|f\|_\beta^\beta}.$$

By taking the expectation w.r.t. $X \sim \mathcal{X}$, we have

$$\begin{aligned} \frac{\mathbb{E}_{X \sim \mathcal{X}} |f(X)g(X)|}{\|f\|_\alpha \|g\|_\beta} &\leq \frac{\mathbb{E}_{X \sim \mathcal{X}} |f(X)|^\alpha}{\alpha \|f\|_\alpha^\alpha} + \frac{\mathbb{E}_{X \sim \mathcal{X}} |f(X)|^\beta}{\beta \|f\|_\beta^\beta} \\ &= \frac{\|f\|_\alpha^\alpha}{\alpha \|f\|_\alpha^\alpha} + \frac{\|f\|_\beta^\beta}{\beta \|f\|_\beta^\beta} \\ &= \frac{1}{\alpha} + \frac{1}{\beta} \\ &= 1. \end{aligned}$$

A.5. HÖLDER's Inequality

This concludes the proof since

$$\frac{\mathbb{E}_{X \sim \mathcal{X}} |f(X)g(X)|}{\|f\|_\alpha \|g\|_\beta} \leq 1 \iff \mathbb{E}_{X \sim \mathcal{X}} |f(X)g(X)| \leq \|f\|_\alpha \|g\|_\beta.$$

■

Abstract. In machine learning, a model is learned from data to solve a task automatically. In the supervised classification setting, the model aims to predict the label associated with an input. The model is learned using a limited number of examples, each consisting of an input and its associated label. However, the model's performance on the examples, computed by the empirical risk, does not necessarily reflect the performance on the task, which is represented by the true risk. Moreover, since it is not computable, the true risk is upper-bounded by a generalization bound that mainly depends on two quantities: the empirical risk and a complexity measure. One way to learn a model is to minimize a bound by a type of algorithm called self-bounding. PAC-Bayesian bounds are well suited to the derivation of this type of algorithm. In this context, the first contribution consists in developing self-bounding algorithms that minimize PAC-Bayesian bounds to learn majority votes. If these bounds are well adapted to majority votes, their use for other models becomes less natural. To overcome this difficulty, a second contribution focuses on the disintegrated PAC-Bayesian bounds that are natural for more general models. In this framework, we provide the first empirical study of these bounds. In a third contribution, we derive bounds that allow us to incorporate complexity measures defined by the user.

Keywords. Machine Learning, Generalization, PAC-Bayesian Bound, Disintegrated PAC-Bayesian Bound, Self-Bounding Algorithm, Majority Vote, Neural Network, Complexity Measure.

Résumé. En apprentissage automatique, un modèle est appris à partir de données pour résoudre une tâche de manière automatique. Dans le cadre de la classification supervisée, le modèle vise à prédire la classe associée à une entrée. Le modèle est appris à l'aide d'un nombre limité d'exemples, chacun étant constitué d'une entrée et de sa classe associée. Cependant, la performance du modèle sur les exemples, calculée par le risque empirique, ne reflète pas nécessairement la performance sur la tâche qui est représentée par le risque réel. De plus, n'étant pas calculable, le risque réel est majoré pour obtenir une borne en généralisation qui dépend principalement de deux quantités : le risque empirique et une mesure de complexité. Une façon d'apprendre un modèle est de minimiser une borne par un type d'algorithme appelé auto-certié (ou auto-limitatif). Les bornes PAC-Bayésiennes sont bien adaptées à la dérivation de ce type d'algorithmes. Dans ce contexte, la première contribution consiste à développer des algorithmes auto-certiés qui minimisent des bornes PAC-Bayésiennes pour apprendre des votes de majorité. Si ces bornes sont bien adaptées aux votes de majorité, leur utilisation pour d'autres modèles devient moins naturelle. Pour pallier cette difficulté, une seconde contribution se concentre sur les bornes PAC-Bayésiennes désintégrées qui sont naturelles pour des modèles plus généraux. Dans ce cadre, nous apportons la première étude empirique de ces bornes. Dans une troisième contribution, nous dérivons des bornes permettant d'incorporer des mesures de complexité pouvant être définies par l'utilisateur.

Mot-clés. Apprentissage Automatique, Généralisation, Borne PAC-Bayésienne, Borne PAC-Bayésienne Désintégrée, Algorithme Auto-certié, Algorithme Auto-limitatif, Vote de Majorité, Réseau de Neurones, Mesure de Complexité.