



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

**On the Interplays between Generalisation and
Optimisation: a PAC-Bayes Approach**

ACKNOWLEDGEMENTS

TODO

Quote 1

Author

Quote 2

AUTHOR

CONTENTS

Contents	5
List of Theorems	8
List of Notations	9
List of Publications	11
Conference article	11
Journal article	11
Research Report	11
1 PAC-Bayes Learning, a field of many paradigms	13
2 PAC-Bayesian Bounds for Martingales and Heavy-Tailed losses	17
3 Mitigating Initialisation Impact by Real-Time Control: Online PAC-Bayes Learning	19
4 Mitigating Initialisation Impact through Flat Minima: Fast Rates for Small Gradients	21
4.1 Introduction	22
5 Wasserstein PAC-Bayes Learning: Exploiting Optimisation Guarantees to Explain Generalisation	23
5.1 Introduction	24
A Some Mathematical Tools	25
A.1 JENSEN's Inequality	25
A.2 MARKOV's Inequality	25
A.3 Ville's Inequality	26
B Additional Background	27
B.1 Online Learning	27
B.2 Wasserstein distances and optimisation on distributions spaces	27
B.3 Differential privacy	27
B.4 Log-Sobolev and Poincaré Inequalities	27

LIST OF ALGORITHMS

Chapter 1

Chapter 2

Chapter 3

Chapter 5

LIST OF THEOREMS

Chapter 1

Chapter 2

Chapter 3

Chapter 4

Chapter 5

Appendix

A.1.1 Theorem (JENSEN's Inequality)	25
A.2.1 Theorem (MARKOV's Inequality)	25
A.3.1 Lemma (Ville's maximal inequality for supermartingales)	26

LIST OF NOTATIONS

General

a	A scalar (integer or real)
\mathbb{R}	The set of real numbers
\mathbb{R}^d	The euclidean set of dimension d
$\ \cdot\ $	a norm of an euclidean set
$\text{dist}(\cdot, \cdot)$	A distance on a Polish space.
\mathbb{N}	The set of natural numbers

Statistical Learning Theory

\mathcal{Z}	Data space. In supervised learning, $\mathcal{Z} = (\mathcal{X}, \mathcal{Y})$ with \mathcal{X}, \mathcal{Y} input and label spaces
\mathbf{z}	A datum of \mathcal{Z} , in supervised learning $\mathbf{z} = (\mathbf{x}, y)$ with \mathbf{x} input and y label
\mathcal{S}	Learning sample $\mathcal{S} = \{\mathbf{z}_i\}_{i \geq 1}$
$\mathcal{D}_{\mathcal{S}}$	Distribution of \mathcal{S}
\mathcal{S}_m	Restriction of \mathcal{S} to its m first data $\mathcal{S}_m = \{\mathbf{z}_i\}_{i=1 \dots m}$
\mathcal{D}	For <i>i.i.d.</i> \mathcal{S} , distribution of a single datum on \mathcal{Z}
\mathcal{D}^m	For <i>i.i.d.</i> \mathcal{S} , distribution of \mathcal{S}_m
\mathcal{T}	For <i>i.i.d.</i> \mathcal{S} , Test set drawn from \mathcal{D}
\mathcal{H}	The set of hypotheses
h	A hypothesis $h \in \mathcal{H}$
ℓ	Loss function $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}$

Probability Theory

$\mathbb{E}_{X \sim \nu} [\cdot]$	The expectation <i>w.r.t.</i> the random variable $X \sim \nu$
$\mathbb{P}_{X \sim \nu} [\cdot]$	The probability <i>w.r.t.</i> the random variable $X \sim \nu$
$\mathbb{1} [a]$	Indicator function; returns 1 if a is true and 0 otherwise
$(\mathcal{F}_i)_{i \geq 1}$	Filtration adapted to \mathcal{S}
$\mathbb{E}_i[\cdot]$	Conditional expectation <i>w.r.t.</i> \mathcal{F}_i , i.e. $\mathbb{E}_i[\cdot] := \mathbb{E}[\cdot \mid \mathcal{F}_i]$
$\mathcal{N}(\mu, \Sigma)$	Gaussian distribution on \mathbb{R}^d with mean μ and covariance matrix Σ

PAC-Bayes framework

$\mathcal{M}(\mathcal{H})$	Set of Probability densities <i>w.r.t.</i> the reference measure on \mathcal{H}
Q	Posterior distribution $Q \in \mathcal{M}(\mathcal{H})$ on \mathcal{H}
P	Prior distribution $P \in \mathcal{M}(\mathcal{H})$ on \mathcal{H}
$KL(Q \parallel P)$	Kullback-Leibler (KL) divergence between Q and P
$D_\alpha(Q \parallel P)$	Rényi Divergence between Q and P
$R_{\mathcal{D}}(h)$	Population Risk of $h \in \mathcal{H}$ <i>w.r.t.</i> \mathcal{D} , i.e. $R_{\mathcal{D}}(h) := \mathbb{E}_{\mathbf{z} \sim \mathcal{D}} [\ell(h, \mathbf{z})]$
$\hat{R}_{\mathcal{S}_m}(h)$	Empirical Risk on \mathcal{S}_m , i.e. $\hat{R}_{\mathcal{S}_m}(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, \mathbf{z}_i)$
$R_{\mathcal{D}}(Q)$	Expected population risk <i>w.r.t.</i> Q , i.e. $R_{\mathcal{D}}(Q) := \mathbb{E}_{h \sim Q} [R_{\mathcal{D}}(h)]$
$\hat{R}_{\mathcal{S}_m}(Q)$	Expected empirical risk <i>w.r.t.</i> Q , $\hat{R}_{\mathcal{S}_m}(Q) := \mathbb{E}_{h \sim Q} [\hat{R}_{\mathcal{S}_m}(h)]$

Optimal transport

W_1	The 1-Wasserstein distance
W_2	The 2-Wasserstein distance
$\Gamma(Q, P)$	Set of all coupling distribution on \mathcal{H}^2 whose marginals are Q and P .

LIST OF PUBLICATIONS

Conference article

PAUL VIALARD, MAXIME HADDOUCHE, Umut SIMSEKLI, and BENJAMIN GUEDJ. Learning via Wasserstein-Based High Probability Generalisation Bounds. *Advances in Neural Information Processing Systems (NeurIPS)*. (2023).

MAXIME HADDOUCHE and BENJAMIN GUEDJ. Online PAC-Bayes Learning. *Advances in Neural Information Processing Systems (NeurIPS)*. (2022).

Journal article

MAXIME HADDOUCHE and BENJAMIN GUEDJ. PAC-Bayes Generalisation Bounds for Heavy-Tailed Losses through Supermartingales. *Transactions on Machine Learning Research*. (2023).

Research Report

MAXIME HADDOUCHE and BENJAMIN GUEDJ. Wasserstein PAC-Bayes Learning: A Bridge Between Generalisation and Optimisation. *arXiv*. abs/2304.07048. (2023).

MAXIME HADDOUCHE, BENJAMIN GUEDJ, and OLIVIER WINTENBERGER. Optimistic Dynamic Regret Bounds. (2023). *arXiv*: 2301.07530 [cs.LG].

PIERRE JOBIC, MAXIME HADDOUCHE, and BENJAMIN GUEDJ. Federated Learning with Nonvacuous Generalisation Bounds. (2023). *arXiv*: 2310.11203 [cs.LG].

MAXIME HADDOUCHE, BENJAMIN GUEDJ, and JOHN SHAWE-TAYLOR. Upper and Lower Bounds on the Performance of Kernel PCA. (2020). *arXiv*: 2012.10369 [cs.LG].

PAC-BAYES LEARNING, A FIELD OF MANY PARADIGMS

1

This manuscript tackles the notion of *generalisation* a notion built upon the general notion of *learning*. For a brief moment, let's take the luxury of forgetting about machines and concentrate on learning at its most human. First and foremost, a human being (or learner) is structured around experiences, either lived or passed on by others. The learner then benefit from these experiences in various ways, for instance, by considering a mediated experience to be true (fire burns) and acting on this assumption, whereas reiteration or denial of this same information may be symptoms of zero truth value. These scenarios can just as easily appear for a lived experience (the question of hallucinations). This first dichotomy in information processing is intrinsically linked to a clearly stated question: does fire burn? Can I trust my senses or have I hallucinated? In these cases, learning has taken place by subjecting the experience to its truth value in relation to a simple question (in this case with two outcomes). This vision can easily be extended to a multiple (and finite) tree of possibilities. Indeed, we can extend the burning question as follows: what is the intensity of the burn as a function of the temperature of the fire? We can then establish a multitude of answers representing various degrees of burn.

However, many questions cannot be reduced to a finite number of possibilities. For example, what is fire? To answer this question, it is nevertheless possible to exploit multiple facets of experience (wood fire, twig, rock) to propose that fire is the chemical reaction of oxygen in the air with a combustible material, with a supply of energy serving as the trigger.

Then, a legitimate question is: why has mankind understood the nature of fire? This fundamental understanding emerged from practical considerations: how can we stop being cold? Can we eat meat other than raw to reduce the risk of illness? It then takes multiple interactions with the environment to generate experiences and then learn from them to gradually respond to a complex need (how to make a fire to keep yourself warm?).

Thus, through this preliminary analysis, we have found several premises of understanding human learning.

- How is learning formalised structurally? The learner must base the experience on simple questions to acquire primary certainties. These latter acquired, it is possible to reach complex questions by interweaving more and more elementary considerations.

- Where does the need to learn come from? From a practical point of view, the emergence of these complex questions often arises from a relationship between the being and its environment, making it possible to develop contextual objectives. The learner then gradually becomes capable of responding to complex needs through a succession of simple actions.

Detail broadly what generalisation is, to what kind of structures it is applied (neural nets or linear classifier eg). Details on the other hand what optimisation is doing (ERM eg) and explain that interestingly in various methods, reaching minimisers of empirical objectives is enough to ensure a good generalisation ability. From this, discuss about the current limitations of generalisation: either not going so often beyond light-tailed assumptions or noticing that the interplays between generalisation (statistical arguments) and optimisation (geometric ones) remains uncharted for a vast range of cases.

Vision: after generic paragraphs on generalisation and optimisation, do a broader paragraph on PAC-Bayes and details the problem of existing PAC-Bayes approach:

Says that PAC-Bayes spontaneously offer a clear link from generalisation to optimisation by providing new learning algorithms: this implicitly suggests assumptions on the loss (eg convex) or on the regulariser (KL between Gaussians to get a strongly convex function) to make sure the minimisation goes well and thus build a bridge with optimisation.

TODO look if there are links from optimisation to PAC-Bayes (must have been some with Dziugaite, Neu with SGD).

Here, we are studying the interplays on both directions. First, we take the opposite perspective and, starting from optimisation benefits/perspectives, we want to understand generalisation, to do so we have several routes within PAC-Bayes. Second, we investigate deeper on the influence of generalisation bounds to derive novel learning algorithms

Thinking the role of the prior in PAC-Bayes: in a similar manner than initialisation/goal to attain in optimisation: if we target a data-free posterior (eg Gibbs distribution) then ok: we target the learning objective. Otherwise, it is common to compare to a random initialisation point of a learning procedure: meaningless. Answer: Online PAC-Bayes which allows, among other, to make the prior evolve through time. (note that there is also either the data-dependent prior: its bad and the differential privacy approach, which is nice)

Switching from statistical to geometric assumptions on the loss. Most of the bounds holds for data-free light-tailed losses: do not necessarily fit the reality of losses involved in optimisation, often unbounded, and either convex, gradient Lipschitz or smooth. Answer: PAC-Bayes for heavy-tailed martingales and flat minima.

Can the convergence properties of optimisation procedure play a role in generalisation? Direct answer: Wasserstein PAC-Bayes, Indirect one: flat minima.

Can we derive generalisation-based learning algorithms beyond Gaussian or Gibbs distributions? Answer: Yes as a byproduct: both in Online PAC-Bayes, Flat Minima and Paper with Paul

Precise the structure of the document: see it as a natural flow where one question implies another one:

Light-tailed losses? \rightarrow supermartingales!

But then : what should we do of the prior? \rightarrow if you see it as an initialisation: Online PAC-Bayes with novel online algorithms or Flat Minima to attenuate the impact of the prior through fast rates.

What if it should be the optimisation goal? \rightarrow Wasserstein PAC-Bayes

CHALLENGE HERE: being very rigorous on the lit review.

PAC-BAYESIAN BOUNDS FOR MARTINGALES AND HEAVY-TAILED LOSSES

This chapter is based on the following paper

MAXIME HADDOUCHE and BENJAMIN GUEDJ. PAC-Bayes Generalisation Bounds for Heavy-Tailed Losses through Supermartingales. *Transactions on Machine Learning Research*. (2023)

Abstract

TODO: put general bounds for martingales and batch learning as corollary

MITIGATING INITIALISATION IMPACT BY REAL-TIME CONTROL: ONLINE PAC-BAYES LEARNING

This chapter is based on the following papers

MAXIME HADDOUCHE and BENJAMIN GUEDJ. Online PAC-Bayes Learning. *Advances in Neural Information Processing Systems (NeurIPS)*. (2022)

MAXIME HADDOUCHE and BENJAMIN GUEDJ. PAC-Bayes Generalisation Bounds for Heavy-Tailed Losses through Supermartingales. *Transactions on Machine Learning Research*. (2023)

PAUL VIALARD, MAXIME HADDOUCHE, UMUT SIMSEKLI, and BENJAMIN GUEDJ. Learning via Wasserstein-Based High Probability Generalisation Bounds. *Advances in Neural Information Processing Systems (NeurIPS)*. (2023)

Contents

4.1	Introduction	22
-----	------------------------	----

Abstract

Put OPB here. Precise in the intro that the martingale bounds allow to go beyond batch learning but that this has never been made for OL. Put the supermartingale OPB bound in a supplementary section and the Online WPB bound after the main results of OPB to reach heavy-tailed losses.

MITIGATING INITIALISATION IMPACT THROUGH FLAT MINIMA: FAST RATES FOR SMALL GRADIENTS

This chapter is based on the following paper

TODO

Contents

5.1	Introduction	24
-----	------------------------	----

Abstract

This is the PLS paper, precise that the supermartingales bounds are richer than simply recovering classical batch guarantees: we can incorporate gradient norms, which explains generalisation when a flat minima is reached.

4.1 Introduction

WASSERSTEIN PAC-BAYES LEARNING: EXPLOITING OPTIMISATION GUARANTEES TO EXPLAIN GENERALISATION

This chapter is based on the following papers

MAXIME HADDOUCHE and BENJAMIN GUEDJ. Wasserstein PAC-Bayes Learning: A Bridge Between Generalisation and Optimisation. *arXiv*. abs/2304.07048. (2023)
PAUL VIALARD, MAXIME HADDOUCHE, Umut SIMSEKLI, and BENJAMIN GUEDJ. Learning via Wasserstein-Based High Probability Generalisation Bounds. *Advances in Neural Information Processing Systems (NeurIPS)*. (2023)

Contents

A.1	JENSEN's Inequality	25
A.2	MARKOV's Inequality	25
A.3	Ville's Inequality	26

Abstract

Put WPB here, precise that, when the prior is seen as the learning goal, it is possible for a certain optimisation algorithm to directly incorporate sound geometric optimisation guarantee into a generalisation bound, trading the hope to reach a flat minima with a sound convergence guarantees. However, this comes at the cost of the explicit impact of the dimension. Also put the paper with Paul(batch bounds) as a supplementary content.

5.1 Introduction

SOME MATHEMATICAL TOOLS

A.1 Jensen's Inequality

Theorem A.1.1 (JENSEN's Inequality). Let X a random variable following a probability distribution ν with f a real-valued measurable convex function, we have

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

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

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

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

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

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

A.2 Markov's Inequality

Theorem A.2.1 (MARKOV's Inequality). Let X a non-negative random variable following a probability distribution ν and $\delta > 0$, we have

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

Proof. First of all, remark that we have the following inequality for any X

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

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

$$\mathbb{E}_{X \sim \nu} [\delta \mathbb{1}[X \geq \delta]] \leq \mathbb{E}_{X \sim \nu} [X].$$

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

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

which is the desired result. ■

A.3 Ville's Inequality

Lemma A.3.1 (Ville's maximal inequality for supermartingales). Let (\mathcal{F}_t) be a filtration and (Z_t) a non-negative super-martingale satisfying $Z_0 = 1$ a.s. If Z_t is adapted to \mathcal{F}_t and $\mathbb{E}[Z_t | \mathcal{F}_{t-1}] \leq Z_{t-1}$ a.s., $t \geq 1$, then, for any $0 < \delta < 1$, it holds

$$\mathbb{P}(\exists T \geq 1 : Z_T > \delta^{-1}) \leq \delta.$$

Proof. We apply the optional stopping theorem (DURRETT, 2019, Thm 4.8.4) with Markov's inequality defining the stopping time $i = \inf\{t \geq 1 : Z_t > \delta^{-1}\}$ so that

$$\mathbb{P}(\exists t \geq 1 : Z_t > \delta^{-1}) = \mathbb{P}(Z_i > \delta^{-1}) \leq \mathbb{E}[Z_i] \delta \leq \mathbb{E}[Z_0] \delta \leq \delta.$$
■

ADDITIONAL BACKGROUND

B

B.1 Online Learning

B.2 Wasserstein distances and optimisation on distributions spaces

B.3 Differential privacy

B.4 Log-Sobolev and Poincaré Inequalities

REFERENCES

RICK DURRETT. Probability: theory and examples. Vol. 49. *Cambridge university press*. (2019)

—— Cited on page 26.

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.