



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

<b>Contents</b>	<b>5</b>
<b>List of Figures</b>	<b>7</b>
<b>List of Theorems</b>	<b>10</b>
<b>List of Notations</b>	<b>11</b>
<b>Preamble: Generalisation and Optimisation in Machine Learning</b>	<b>13</b>
<b>List of Publications</b>	<b>15</b>
Conference article . . . . .	15
Journal article . . . . .	15
Research Report . . . . .	15
 <b>I Background</b>	 <b>17</b>
<b>1 An introduction to Statistical Learning Theory and PAC-Bayes</b>	<b>19</b>
1.1 A brief survey of generalisation bounds . . . . .	20
1.2 PAC-Bayes learning . . . . .	20
 <b>2 A Brief Reminder on Optimisation for Batch and Online Learning</b>	 <b>21</b>
2.1 Convergence guarantees for optimisation on predictor spaces . . . . .	22
2.2 Wasserstein distances and optimisation on measure spaces . . . . .	22
2.3 Optimisation in Online Learning . . . . .	22
 <b>3 Additionnal Tools to Link Generalisation and Optimisation</b>	 <b>23</b>
3.1 Differential privacy . . . . .	23
3.2 Log-Sobolev inequalities . . . . .	23
 <b>II Generalisation bounds for Martingales and Online Learning allowing Heavy-Tailed Losses</b>	 <b>25</b>
 <b>4 PAC-Bayesian Bounds for Martingales and Heavy-Tailed losses</b>	 <b>27</b>

<b>5</b>	<b>Online PAC-Bayes Learning for Bounded Losses and Beyond</b>	<b>29</b>
<b>III</b>	<b>Towards A Better Understanding of Generalisation through Optimisation</b>	<b>31</b>
<b>6</b>	<b>Understanding Generalisation through Flat Minimas</b>	<b>33</b>
6.1	Introduction . . . . .	34
<b>7</b>	<b>Wasserstein PAC-Bayes Learning: Exploiting Optimisation Guarantees to Explain Generalisation</b>	<b>35</b>
7.1	Introduction . . . . .	36
<b>IV</b>	<b>Conclusion and Perspectives</b>	<b>37</b>
<b>V</b>	<b>Appendix</b>	<b>39</b>
<b>A</b>	<b>Some Mathematical Tools</b>	<b>41</b>
A.1	JENSEN's Inequality . . . . .	41
A.2	MARKOV's Inequality . . . . .	41
A.3	2nd Order MARKOV's Inequality . . . . .	42
A.4	CHEBYSHEV-CANTELLI Inequality . . . . .	43
A.5	HÖLDER's Inequality . . . . .	43

# LIST OF FIGURES

**Preamble**

**Chapter 1**

**Chapter 2**

**Chapter 3**

**Chapter 4**

**Chapter 5**

**Chapter 6**

**Chapter 7**





# LIST OF ALGORITHMS

**Chapter 1**

**Chapter 2**

**Chapter 3**

**Chapter 4**

**Chapter 5**

**Chapter 7**

# LIST OF THEOREMS

## Chapter 1

## Chapter 2

## Chapter 3

## Chapter 4

## Chapter 5

## Chapter 6

## Chapter 7

## Appendix

A.1.1 Theorem (JENSEN's Inequality) . . . . .	41
A.2.1 Theorem (MARKOV's Inequality) . . . . .	41
A.3.1 Theorem (2nd Order MARKOV's Inequality) . . . . .	42
A.4.1 Theorem (CHEBYSHEV-CANTELLI Inequality) . . . . .	43
A.5.1 Lemma (YOUNG's Inequality) . . . . .	43
A.5.1 Theorem (HÖLDER's Inequality) . . . . .	44

# 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 $\mathbb{H}$
$\rho$	Posterior distribution $\rho \in \mathbb{M}(\mathbb{H})$ on $\mathbb{H}$
$\pi$	Prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$ on $\mathbb{H}$
$\text{KL}(\rho \parallel \pi)$	Kullback-Leibler (KL) divergence between $\rho$ and $\pi$
$D_\alpha(\rho \parallel \pi)$	Rényi Divergence between $\rho$ and $\pi$
$\text{Uni}(\mathbb{A})$	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 $\text{MV}_\rho(\cdot)$

# PREAMBLE: GENERALISATION AND OPTIMISATION IN MACHINE LEARNING

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 catoni) 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.

# 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].





## PART I

# **Background**



# AN INTRODUCTION TO STATISTICAL LEARNING THEORY AND PAC-BAYES

## Contents

1.1	A brief survey of generalisation bounds . . . . .	20
1.2	PAC-Bayes learning . . . . .	20

## Abstract

This chapter provides a brief introduction of statistical learning theory and various modern kind of generalisation bounds (go from uniform convergence bounds to algorithmic stability, search for other kinds of non PAC-Bayes bounds). Need to recall the historical and modern shapes of generalisation bounds in ML.

## **1.1 A brief survey of generalisation bounds**

## **1.2 PAC-Bayes learning**

# A BRIEF REMINDER ON OPTIMISATION FOR BATCH AND ONLINE LEARNING

## Contents

---

2.1	Convergence guarantees for optimisation on predictor spaces . . . . .	<b>22</b>
2.2	Wasserstein distances and optimisation on measure spaces . . . . .	<b>22</b>
2.3	Optimisation in Online Learning . . . . .	<b>22</b>

---

## Abstract

Detail optimisation for GD, SGD and variants, list a lot of convergence guarantees under various assumptions (convexity smoothness etc). On the measure spaces part, detail the required OT background to introduce Wasserstein distances. Detail what is online learning.

- 2.1 Convergence guarantees for optimisation on predictor spaces**
- 2.2 Wasserstein distances and optimisation on measure spaces**
- 2.3 Optimisation in Online Learning**

# ADDITIONNAL TOOLS TO LINK GENERALISATION AND OPTIMISATION

# 3

## Contents

3.1	Differential privacy . . . . .	23
3.2	Log-Sobolev inequalities . . . . .	23

## Abstract

Put here additionnal background on differential privacy and log-Sobolev inequalities.

## 3.1 Differential privacy

## 3.2 Log-Sobolev inequalities





## PART II

# **Generalisation bounds for Martingales and Online Learning allowing Heavy-Tailed Losses**



# 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



# ONLINE PAC-BAYES LEARNING FOR BOUNDED LOSSES AND BEYOND

# 5

**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)

## **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.



## PART III

# **Towards A Better Understanding of Generalisation through Optimisation**





# UNDESTANDING GENERALISATION THROUGH FLAT MINIMAS

This chapter is based on the following paper

TODO

## Contents

---

6.1	Introduction . . . . .	34
-----	------------------------	----

---

## 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.

## 6.1 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

7.1	Introduction . . . . .	36
-----	------------------------	----

## Abstract

Put WPB here, precise that beyond the somewhat vague understanding of generalisation through flat minima, it is possible, for a certain optimisation algorithm to directly incorporate a 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

## 7.1 Introduction

## 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.