

PAC-Bayes bounds: an introduction

Pierre Alquier



Post-Bayes seminar series
Chapter 3: PAC-Bayes
September 23, 2025

Welcome to the post-Bayes seminar series !

- At a glance/website : [▶ Link](https://tinyurl.com/postBayesWebsite) <https://tinyurl.com/postBayesWebsite>
- Where to subscribe to mailing list : [▶ Link](#)
- Where to subscribe to calendard : [▶ Link](#)
- Where to attend the seminars : [▶ Link](#)
- Where recorded seminars are stored : [▶ Link](#)

Please share widely !

Welcome to the post-Bayes seminar series !

During talk :

- Use Q/A function in zoom
- Other questions can be upvoted
- We will try to monitor questions and ask relevant ones in natural breaks

After talk :

- Raise your hand in zoom
- We will do our best to decide who gets to ask a question fairly
- We will do our best to resolve remaining questions in Q / A function

Structure of Chapter 3 : PAC-Bayes

	Prof. Pierre Alquier	Today PAC-Bayes : an introduction
	Prof. Dan Roy	30/09 – Removing models and assumptions from Bayesian Decision Making
	Prof. Yevgeny Seldin	07/10 Recursive PAC-Bayes
	Prof. Pascal Germain	28/10 PAC-Bayes Hypernetworks
	Prof. Benjamin Guedj	04/11 Rethinking Generalisation : Beyond KL with Geometry and Comparators
	Dr. Badr-Eddine Chérief-Abdellatif	18/11 PAC-Bayes Meets Variational Inference : Theory and Generalizations

Cergy (near Paris)



- Kamelia DAUDEL
- Olga KLOPP
- Marie KRATZ
- Guillaume CHEVILLON
- Roberto RENO
- Mikolaj KASPRZAK
- Jeroen ROMBOUTS
- Vincenzo ESPOSITO VINZI
- Mohamed NDAOUD
- Guillaume LECUE
- Pierre JACOB
- Maria ALLAYIOTI

Singapore



- Jeremy HENG
- Pierre ALQUIER

1 PAC-Bayes bounds : introduction

- Generalization bounds and PAC-Bayes
- Minimization of the PAC-Bayes bound
- References

2 Relevance of PAC-Bayes in the post-Bayes community

- Rates of convergence
- Analysis of generalized posteriors
- Mutual Information bounds : optimizing the prior

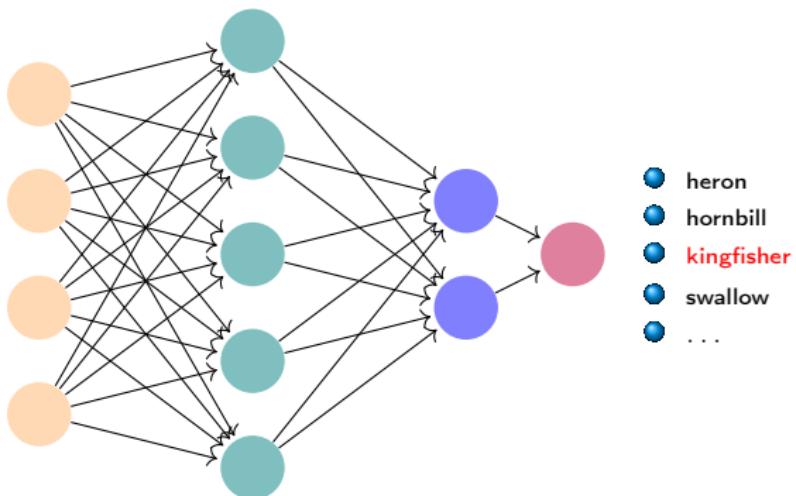
1 PAC-Bayes bounds : introduction

- Generalization bounds and PAC-Bayes
- Minimization of the PAC-Bayes bound
- References

2 Relevance of PAC-Bayes in the post-Bayes community

- Rates of convergence
- Analysis of generalized posteriors
- Mutual Information bounds : optimizing the prior

- Objects $x \in \mathcal{X}$, labels $y \in \mathcal{Y}$.
- Predictor : function $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ indexed by $\theta \in \Theta$.



- Prediction error measured through loss function ℓ :

$$\ell\left(y, f_{\theta}(x)\right).$$

- Risk :

$$R(\theta) := \mathbb{E}_{(X, Y) \sim P} \left[\ell\left(Y, f_{\theta}(X)\right) \right].$$

where P is the probability distribution of pairs object-label we want to learn to classify.

- Objective :

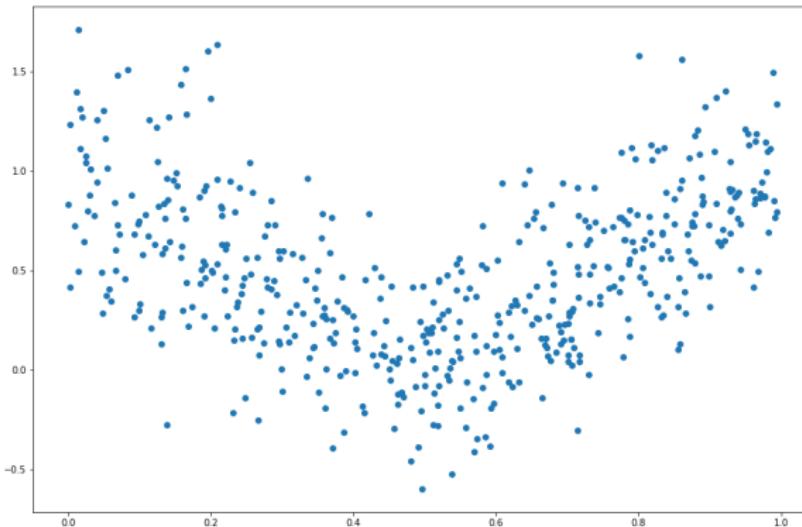
$$R^* = \inf_{\theta \in \Theta} R(\theta).$$

- Data $(X_1, Y_1), \dots, (X_n, Y_n)$ i.i.d. from P . Empirical risk :

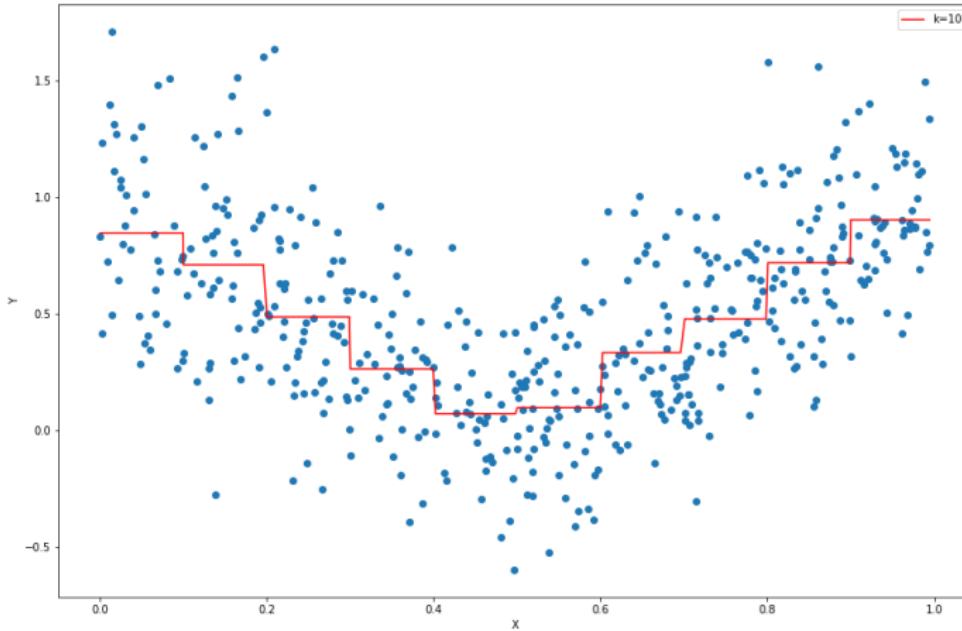
$$R_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell\left(Y_i, f_{\theta}(X_i)\right).$$

Toy example :

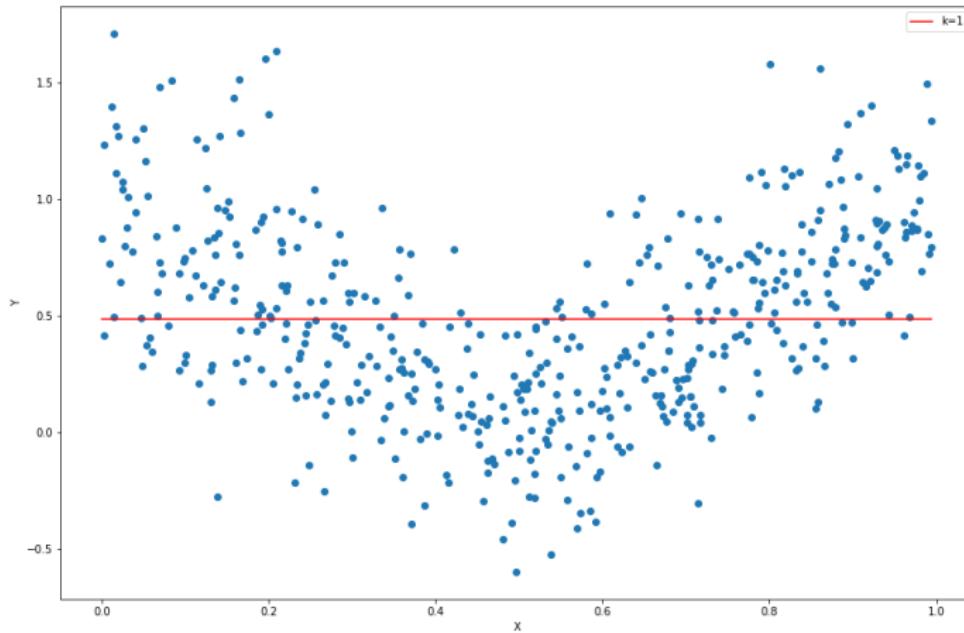
- X uniform on $[0, 1]$,
- $Y = |2X - 1| + \epsilon$.



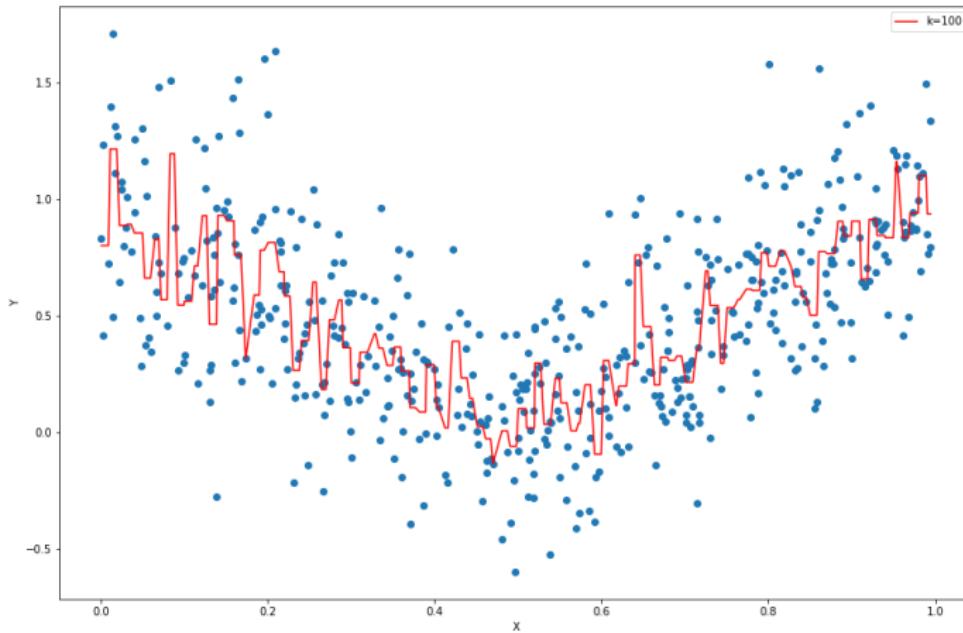
- Prediction by regular histogram with k -bins.
- $\ell(y, f_\theta(x)) = (y - f_\theta(x))^2$.

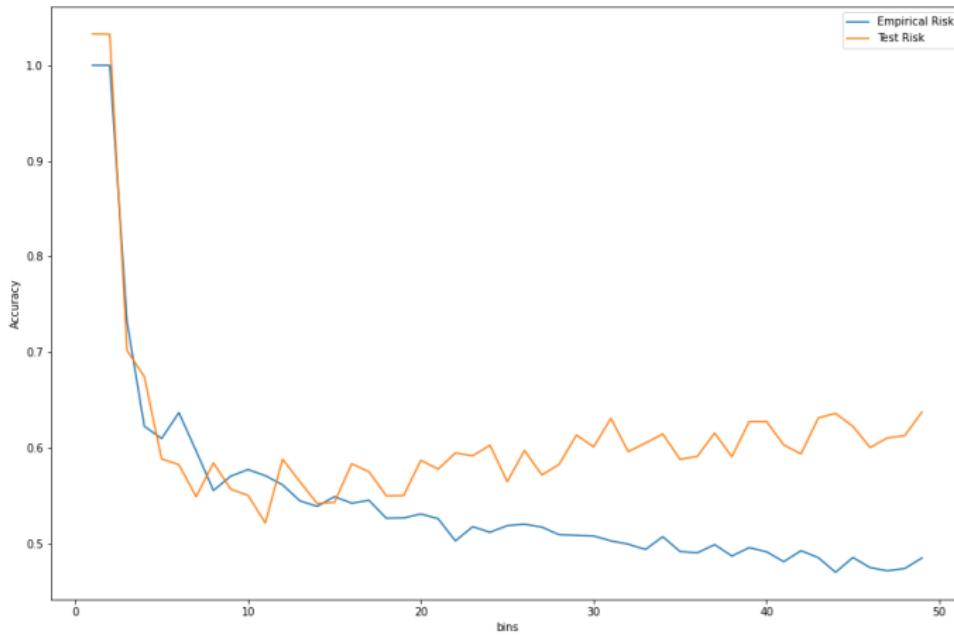


- Prediction by regular histogram with k -bins.
- $\ell(y, f_\theta(x)) = (y - f_\theta(x))^2$.



- Prediction by regular histogram with k -bins.
- $\ell(y, f_\theta(x)) = (y - f_\theta(x))^2$.





Law of large numbers : for a fixed θ ,

$$R_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell\left(Y_i, f_\theta(X_i)\right) \xrightarrow{n \rightarrow \infty} R(\theta).$$

But $\hat{\theta} = \hat{\theta}((X_1, Y_1), \dots, (X_n, Y_n)) = \hat{\theta}(\mathcal{S})$ learnt from data.

Can we quantify $R(\hat{\theta}) - R_n(\hat{\theta})$ when $\hat{\theta}$ is learnt ?

Various approaches :

- Vapnik-Chervonenkis theory,
- algorithmic stability,
- information bounds : MDL, PAC-Bayes, etc.

Assumption for whole lecture

Unless specified otherwise, $0 \leq \ell \leq 1$ and data is i.i.d. from P .

Vapnik-Chervonenkis – classification ($\mathcal{Y} = \{0, 1\}$)

With probability at least $1 - \delta$ on the data, for any $\hat{\theta}$ learnt from the data,

$$R(\hat{\theta}) \leq R_n(\hat{\theta}) + \sqrt{\frac{8d \log\left(\frac{2en}{d}\right) + 8 \log\left(\frac{4}{\delta}\right)}{n}}$$

where d : the VC-dimension of the set of classifiers ($f_\theta, \theta \in \Theta$).

Statistical estimation / ERM etc.

data \longrightarrow estimator

$$(\mathcal{X} \times \mathcal{Y})^n \longrightarrow \Theta$$

$$\mathcal{S} \longleftarrow \hat{\theta} = \hat{\theta}(\mathcal{S})$$

Randomized estimators :

$$(\mathcal{X} \times \mathcal{Y})^n \longrightarrow \mathcal{M}(\Theta) \dashrightarrow \Theta$$

$$\mathcal{S} \longleftarrow \hat{\theta} = \hat{\theta}(\mathcal{S}) \dashrightarrow^{\theta \sim \hat{\rho}} \theta$$

McAllester's PAC-Bayes bound

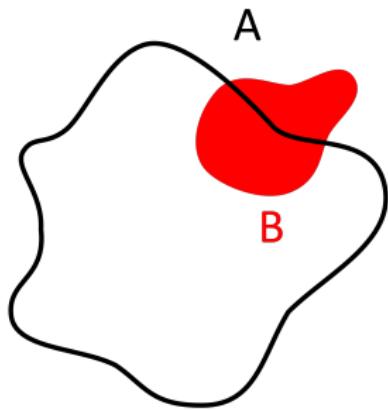
Fix a prior distribution $\pi \in \mathcal{M}(\Theta)$. With probability at least $1 - \delta$ on the data \mathcal{S} , for any probability distribution ρ learnt on the data,

$$\mathbb{E}_{\theta \sim \rho}[R(\theta)] \leq \mathbb{E}_{\theta \sim \rho}[R_n(\theta)] + \sqrt{\frac{\text{KL}(\rho \| \pi) + \log\left(\frac{2\sqrt{n}}{\delta}\right)}{2n}}.$$

$\text{KL}(\rho \| \pi) =$ Kullback-Leibler divergence between ρ and π

- ρ can be learnt on the data, so if we have a randomized estimator $\hat{\rho}$ in mind, we can apply the bound to $\rho = \hat{\rho}$.
- we will see later that the bound is helpful to define good randomized estimators $\hat{\rho}$.

Intuition on KL :



- π uniform on A

$$\pi(\theta) = \frac{1_A(\theta)}{\mathcal{V}(A)}$$

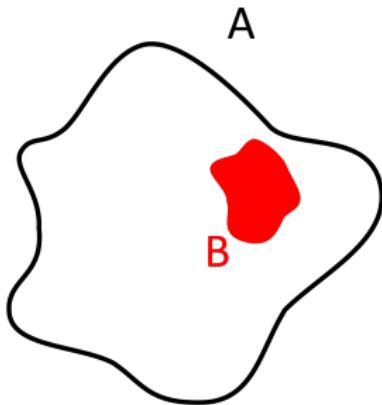
- ρ uniform on B

$$\rho(\theta) = \frac{1_B(\theta)}{\mathcal{V}(B)}$$

$\frac{d\rho}{d\pi}$ not defined here.

$B \not\subseteq A \Rightarrow \text{KL}(\rho\|\pi) = +\infty.$

Intuition on KL :



- π uniform on A

$$\pi(\theta) = \frac{1_A(\theta)}{\mathcal{V}(A)}$$

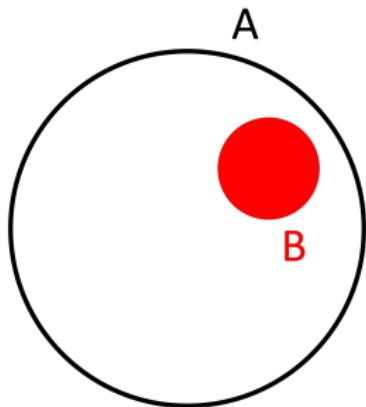
- ρ uniform on B

$$\rho(\theta) = \frac{1_B(\theta)}{\mathcal{V}(B)}$$

$$B \subseteq A \Rightarrow \frac{d\rho}{d\pi}(\theta) = \frac{\mathcal{V}(A)1_B(\theta)}{\mathcal{V}(B)}$$

$$KL(\rho\|\pi) = \mathbb{E}_{\theta \sim \rho} \left[\log \frac{d\rho}{d\pi}(\theta) \right] = \log \frac{\mathcal{V}(A)}{\mathcal{V}(B)}.$$

Intuition on KL :



$B_d(x, r)$ ball centered on x , with radius r in \mathbb{R}^d

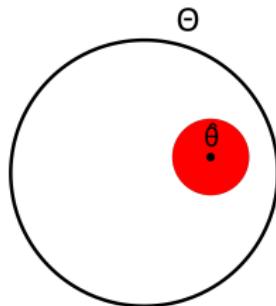
$$\mathcal{V}(B_d(x, r)) = \frac{r^d \pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2} + 1)}$$

- π uniform on $A = B_d(0, C)$
- ρ uniform on $B = B_d(\theta_0, \epsilon)$

$$\text{KL}(\rho \| \pi) = \log \frac{\mathcal{V}(A)}{\mathcal{V}(B)} = d \log \frac{C}{\epsilon}.$$

McAllester's PAC-Bayes bound

$$\mathbb{E}_{\theta \sim \rho}[R(\theta)] \leq \mathbb{E}_{\theta \sim \rho}[R_n(\theta)] + \sqrt{\frac{\text{KL}(\rho \parallel \pi) + \log \left(\frac{2\sqrt{n}}{\delta} \right)}{2n}}.$$



- π uniform on $\Theta = B_d(0, C)$
- $\rho = \hat{\rho}$ uniform on $B_d(\hat{\theta}, \epsilon)$

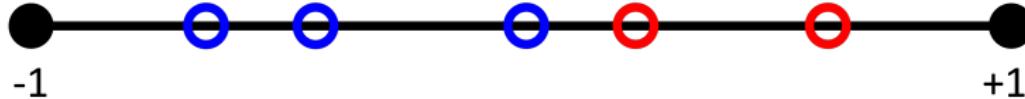
$$\mathbb{E}_{\theta \sim \hat{\rho}}[R(\theta)] \leq \mathbb{E}_{\theta \sim \hat{\rho}}[R_n(\theta)] + \sqrt{\frac{d \log \frac{C}{\epsilon} + \log \left(\frac{2\sqrt{n}}{\delta} \right)}{2n}}.$$

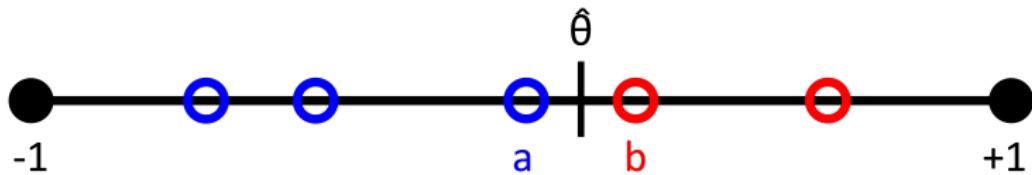
Toy classification example :

- $X_i \in [-1, 1]$,
- classifiers $(f_\theta)_{\theta \in [-1, 1]}$ given by

$$f_\theta(x) = \begin{cases} 0 & \text{if } x \leq \theta \\ 1 & \text{if } x > \theta. \end{cases}$$

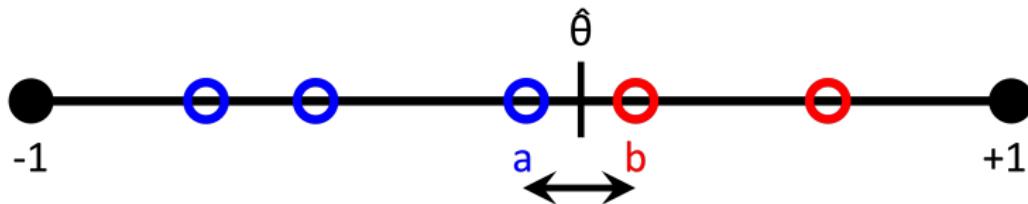
- $Y_i = f_{\theta^*}(X_i)$.





Vapnik-type bound :

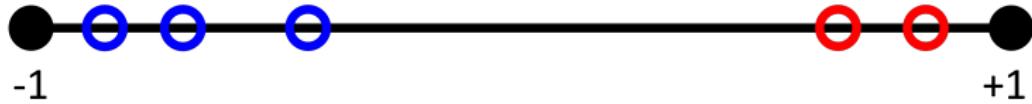
$$R(\hat{\theta}) \leq \sqrt{\frac{8 \log(2en) + 8 \log\left(\frac{4}{\delta}\right)}{n}}$$

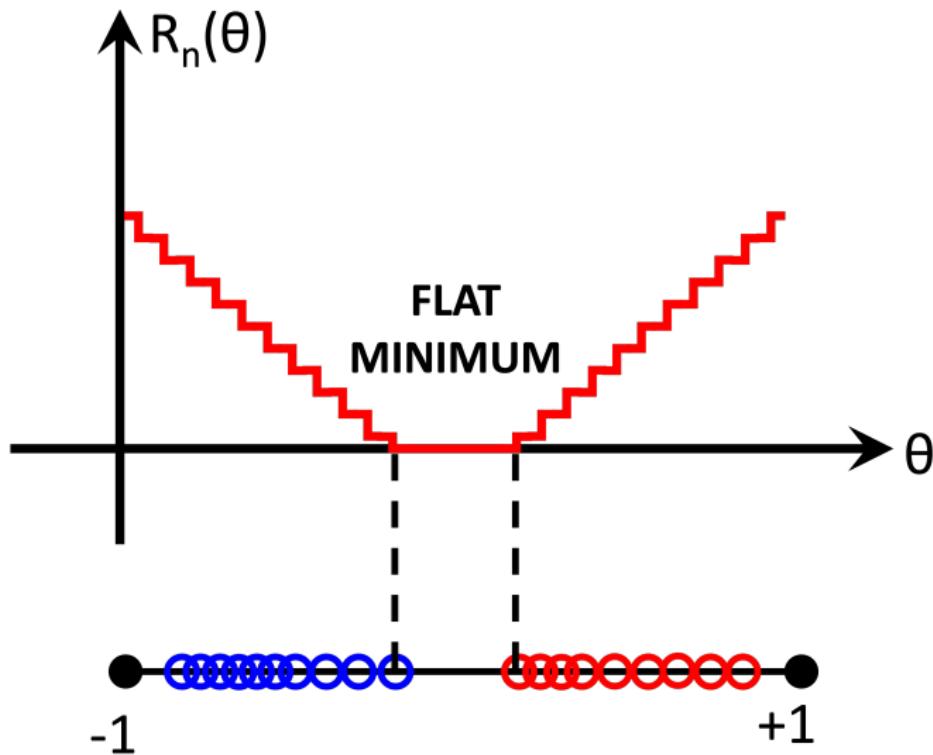


PAC-Bayes :

- π uniform on $[-1, 1]$,
- $\hat{\rho}$ uniform on $[a, b]$.

$$\mathbb{E}_{\theta \sim \hat{\rho}}[R(\theta)] \leq \sqrt{\frac{\log \frac{2}{b-a} + \log \left(\frac{2\sqrt{n}}{\delta} \right)}{2n}}.$$





1 PAC-Bayes bounds : introduction

- Generalization bounds and PAC-Bayes
- Minimization of the PAC-Bayes bound
- References

2 Relevance of PAC-Bayes in the post-Bayes community

- Rates of convergence
- Analysis of generalized posteriors
- Mutual Information bounds : optimizing the prior

McAllester's PAC-Bayes bound

Fix prior $\pi \in \mathcal{M}(\Theta)$. With proba. at least $1 - \delta$, $\forall \rho \in \mathcal{M}(\Theta)$,

$$\mathbb{E}_{\theta \sim \rho}[R(\theta)] \leq \mathbb{E}_{\theta \sim \rho}[R_n(\theta)] + \sqrt{\frac{\text{KL}(\rho \| \pi) + \log\left(\frac{2\sqrt{n}}{\delta}\right)}{2n}}.$$

$$\sqrt{\frac{a}{b}} = \inf_{\lambda > 0} \left\{ \frac{a}{\lambda} + \frac{\lambda}{4b} \right\}.$$

$$\mathbb{E}_{\theta \sim \rho}[R(\theta)] \leq \mathbb{E}_{\theta \sim \rho}[R_n(\theta)]$$

$$+ \inf_{\lambda > 0} \left\{ \frac{\text{KL}(\rho \| \pi) + \log\left(\frac{2\sqrt{n}}{\delta}\right)}{\lambda} + \frac{\lambda}{8n} \right\}.$$

$$\forall \lambda > 0, \mathbb{E}_{\theta \sim \rho}[R(\theta)] \leq \mathbb{E}_{\theta \sim \rho}[R_n(\theta)]$$

$$+ \frac{\text{KL}(\rho \| \pi) + \log\left(\frac{2\sqrt{n}}{\delta}\right)}{\lambda} + \frac{\lambda}{8n}.$$

Definition - Gibbs posterior

$$\hat{\pi}_\lambda(d\theta) = \frac{\exp(-\lambda R_n(\theta))}{\mathbb{E}_{\vartheta \sim \pi}[\exp(-\lambda R_n(\vartheta))]} \pi(d\theta).$$

Theorem

$$\hat{\pi}_\lambda = \arg \min_{\rho \in \mathcal{M}(\Theta)} \left\{ \mathbb{E}_{\theta \sim \rho}[R_n(\theta)] + \frac{\text{KL}(\rho \| \pi)}{\lambda} \right\}.$$

Sampling from $\hat{\pi}_\lambda$ by Monte-Carlo techniques...

Approximate minimization of the PAC-Bayes bound.

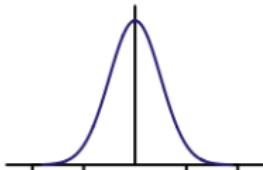
$$\forall \lambda > 0, \mathbb{E}_{\theta \sim \rho}[R(\theta)] \leq \mathbb{E}_{\theta \sim \rho}[R_n(\theta)]$$

$$+ \frac{\text{KL}(\rho \| \pi) + \log \left(\frac{2\sqrt{n}}{\delta} \right)}{\lambda} + \frac{\lambda}{8n}.$$

Alternative approach : optimize ρ in a smaller set $\mathcal{F} \subsetneq \mathcal{M}(\Theta)$.

Definition - variational approximation of Gibbs posterior

$$\tilde{\rho}_\lambda = \arg \min_{\rho \in \mathcal{F}} \left\{ \mathbb{E}_{\theta \sim \rho}[R_n(\theta)] + \frac{\text{KL}(\rho \| \pi)}{\lambda} \right\}.$$



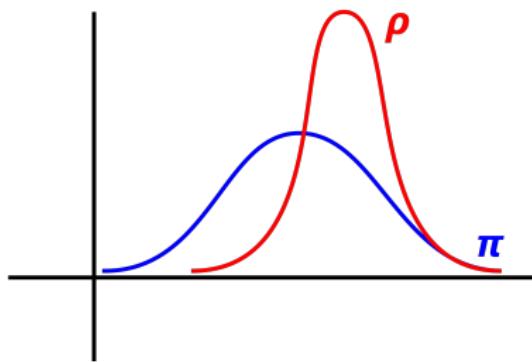
Example : $\rho = \mathcal{N}(\mu, \Sigma)$, optimize (μ, Σ) .

Example : Gaussian prior π , and we optimize a Gaussian posterior ρ :

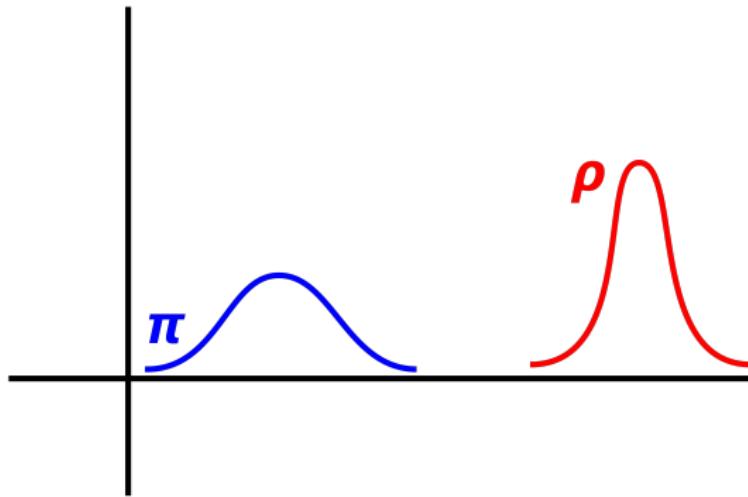
$$\pi = \mathcal{N}(\mu_0, \Sigma_0) \text{ and } \rho = \mathcal{N}(\mu_1, \Sigma_1) \text{ in } \mathbb{R}^d.$$

$$\begin{aligned} \text{KL}(\rho\|\pi) &= \frac{1}{2} \left[\text{tr}(\Sigma_1 \Sigma_0^{-1}) - d \right. \\ &\quad \left. + (\mu_1 - \mu_0)^T \Sigma_0^{-1} (\mu_1 - \mu_0) + \log \frac{\det \Sigma_0}{\det \Sigma_1} \right]. \end{aligned}$$

$\pi = \mathcal{N}(\mu_0, \Sigma_0)$ and $\rho = \mathcal{N}(\mu_1, \Sigma_1)$ in \mathbb{R} .

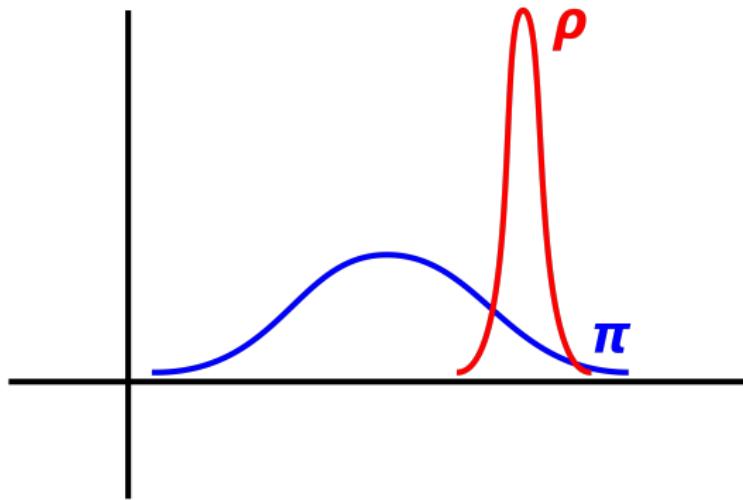


$$\text{KL}(\rho\|\pi) = \frac{1}{2} \left[\frac{\Sigma_1}{\Sigma_0} - 1 + \frac{(\mu_0 - \mu_1)^2}{\Sigma_0} + \log \frac{\Sigma_0}{\Sigma_1} \right].$$



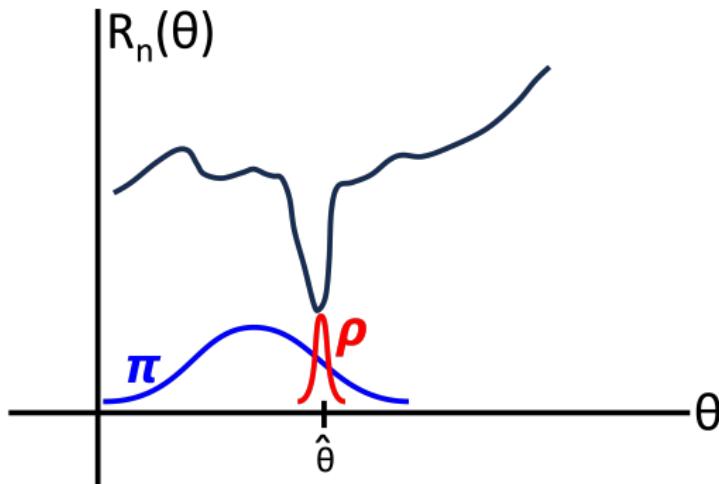
If μ_1 goes far away from μ_0 to ∞ ,

$$\text{KL}(\rho \parallel \pi) \sim \frac{(\mu_0 - \mu_1)^2}{2\Sigma_0} \rightarrow \infty.$$



If $\Sigma_1 \rightarrow 0$,

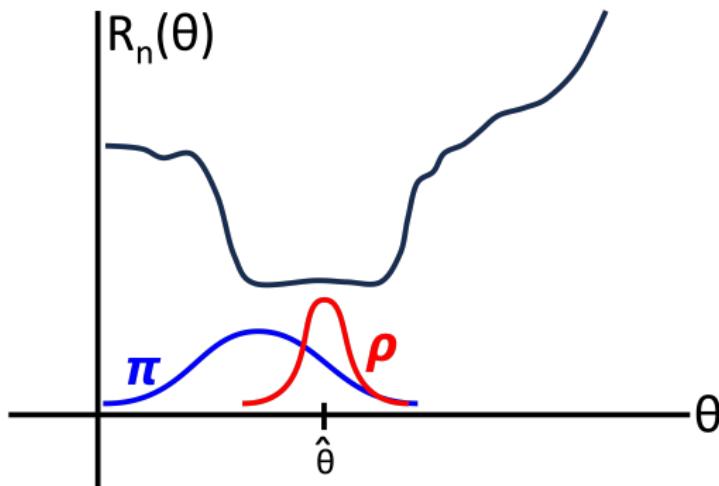
$$\text{KL}(\rho\|\pi) \sim \frac{1}{2} \log \frac{\Sigma_0}{\Sigma_1} \rightarrow \infty.$$



With a sharp minimum, to keep

$$\mathbb{E}_{\theta \sim \mathcal{N}(\hat{\theta}, \Sigma_1)} [R_n(\theta)] \sim R_n(\hat{\theta}),$$

Σ_1 should be small, and thus $\text{KL}(\rho \parallel \pi)$ will be large.



With a flat minimum,

$$\mathbb{E}_{\theta \sim \mathcal{N}(\hat{\theta}, \Sigma_1)} [R_n(\theta)] \sim R_n(\hat{\theta})$$

for Σ_1 “not so small”, thus $\text{KL}(\rho \parallel \pi)$ does not have to be large.

Application : generalization bounds for deep learning.

Train a neural network for classification (0-1 loss).

Vapnik-type bound usually lead to something larger than 1, for example :

$$R(\hat{\theta}) \leq 35.4$$

As $R(\hat{\theta}) = \mathbb{P}(Y \neq f_{\hat{\theta}}(X))$, the bound brings no information (vacuous).

Computing Nonvacuous Generalization Bounds for Deep (Stochastic) Neural Networks with Many More Parameters than Training Data

Gintare Karolina Džingaitė
Department of Engineering
University of Cambridge

Daniel M. Roy
Department of Statistical Sciences
University of Toronto

Abstract

One of the defining properties of deep learning is that models are chosen to have many more parameters than available training data. In light of this capacity for overfitting, it is remarkable that simple algorithms like SGD reliably return solutions with low test error. One roadblock to explaining these phenomena in terms of implicit regularization, structural properties of the solution, and/or easiness of the data is that many learning bounds are quantitatively vacuous when applied to networks learned by SGD in this “deep learning” regime. Logically, in order to explain generalization, we need nonvacuous bounds. We return to an idea by Langford and Caruana (2001), who used PAC-Bayes bounds to compute nonvacuous numerical bounds on generalization error for *stochastic* two-layer two-hidden-unit neural networks via a sensitivity analysis. By optimizing the PAC-Bayes bound directly, we are able to extend their approach and obtain nonvacuous generalization bounds for deep stochastic neural network classifiers with millions of parameters trained on only tens of thousands of examples. We connect our findings to recent and old work on flat minima and MDL-based explanations of generalization.

1 INTRODUCTION

By optimizing a PAC-Bayes bound, we show that it is possible to compute nonvacuous numerical bounds on the generalization error of deep stochastic neural networks with millions of parameters, despite the training data sets being one or more orders of magnitude smaller than the number of parameters. To our knowledge, these are the first explicit and nonvacuous numerical bounds computed

for trained neural networks in the modern deep learning regime where the number of network parameters eclipses the number of training examples.

The bounds we compute are dependent, incorporating millions of components optimized numerically to identify a large region in weight space with low average empirical error around the solutions obtained by stochastic gradient descent (SGD). The data dependence is essential; indeed, the VC dimension of neural networks is typically bounded below by the number of parameters, and so one needs as many training data as parameters before (uniform) PAC bounds are nonvacuous, i.e., before the generalization error falls below 1. To put this in concrete terms, on MNIST, having even 72 hidden units in a fully connected first layer yields vacuous PAC bounds.

Evidently, we are operating far from the worst case: observed generalization cannot be explained in terms the regularizing effect of the size of the neural network alone. This is an old observation, and one that attracted considerable theoretical attention two decades ago: Bartlett [Bar97; Bar98] showed that, in large (sigmoidal) neural networks, when the learned weights are small in magnitude, the fat-shattering dimension is more important than the VC dimension for characterizing generalization. In particular, Bartlett established classification error bounds in terms of the empirical margin and the fat-shattering dimension, and then gave fat-shattering bounds for neural networks in terms of the magnitudes of the weights and the depth of the network alone. Improved norm-based bounds were obtained using Rademacher and Gaussian complexity by Bartlett and Mendelson [BM02] and Koltchinskii and Panchenko [KP02].

These norm-based bounds are the foundation of our current understanding of neural network generalization. It is widely accepted that these bounds explain observed generalization, at least “qualitatively” and/or when the weights are explicitly regularized. Indeed, recent work by Neyshabur, Tomioka, and Srebro [NTS14] puts forth

Experiment	T-600	T-1200	T-300 ²	T-600 ²	T-1200 ²	T-600 ³	R-600
Train error	0.001	0.002	0.000	0.000	0.000	0.000	0.007
Test error	0.018	0.018	0.015	0.016	0.015	0.013	0.508
SNN train error	0.028	0.027	0.027	0.028	0.029	0.027	0.112
SNN test error	0.034	0.035	0.034	0.033	0.035	0.032	0.503
PAC-Bayes bound	0.161	0.179	0.170	0.186	0.223	0.201	1.352
KL divergence	5144	5977	5791	6534	8558	7861	201131
# parameters	471k	943k	326k	832k	2384k	1193k	472k
VC dimension	26m	56m	26m	66m	187m	121m	26m

Table 1: Results for experiments on binary class variant of MNIST. SGD is either trained on (T) true labels or (R) random labels. The network architecture is expressed as N^L , indicating L hidden layers with N nodes each. Errors are classification error. The reported VC dimension is the best known upper bound (in millions) for ReLU networks. The SNN error rates are tight upper bounds (see text for details). The PAC-Bayes bounds upper bound the test error with probability 0.965.

Results taken from :



Dzuigaite, G. K. and Roy, D. M. (2017). Computing Nonvacuous Generalization Bounds for Deep (Stochastic) Neural Networks with Many More Parameters than Training Data. *UAI*.

1 PAC-Bayes bounds : introduction

- Generalization bounds and PAC-Bayes
- Minimization of the PAC-Bayes bound
- References

2 Relevance of PAC-Bayes in the post-Bayes community

- Rates of convergence
- Analysis of generalized posteriors
- Mutual Information bounds : optimizing the prior

PAC-Bayesian Model Averaging

David A. McAllester
AT&T Shannon Labs
180 Park Avenue
Florham Park, NJ 07932-0971
dmac@research.att.com

Abstract

PAC-Bayesian learning methods combine the informative priors of Bayesian methods with distribution-free PAC guarantees. Building on earlier methods for PAC-Bayesian model selection, this paper presents a method for PAC-Bayesian model averaging. The method constructs an optimized weighted mixture of competing models that minimizes the error of overfit. Although the main result is aimed for bounded loss, a preliminary analysis for unbounded loss is also given.

1 INTRODUCTION

A PAC-Bayesian approach to machine learning attempts to combine the advantages of both PAC and Bayesian approaches [12, 8]. The Bayesian approach has the advantage of using arbitrary domain knowledge in the form of a Bayesian prior. The PAC approach has the advantage that one can prove guarantees for generalization error based on the truth of the prior. The PAC-Bayesian approach combines the features of the PAC and Bayesian approaches: it bases the bias of the learning algorithm on an arbitrary prior distribution, thus allowing the incorporation of domain knowledge, and yet provides a guarantee on generalization error that is based on the truth of the prior.

PAC-Bayesian approaches are related to structural risk minimization (SRM) [6]. Here we interpret this broadly as describing any learning algorithm optimizing a trade-off between “complexity”, “structure”, “prior probability” of the concepts in mind, and the “goodness of fit”, “description length”, or “likelihood” of the training data. Under this interpretation of SRM, Bayesian algorithms which select a concept of maximum posterior probability (MAP algorithms) are viewed as a kind of SRM algorithm. Various approaches to SRM

Permission to make digital or hard copies of all or part of this work for personal use or classroom use without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
COLT '99 7-9 June Santa Cruz, CA, USA
© 1999 ACM 1-58113-167-4/99/0006-08-00

are compared both theoretically and experimentally by Kearns et al. in [6]. They give experimental evidence that Bayesian and MDL algorithms tend to over-fit in experiments involving noisy Bayesian networks. A PAC-Bayesian approach uses a prior distribution analogous to that used in MAP or MDL but provides a theoretical guarantee against over-fitting independent of the truth of the prior.

Earlier work on PAC-Bayesian algorithms has focused on model selection, either a single model except or a uniformly weighted set of concepts. Here we consider nonuniform model averaging, i.e., selecting a weighted mixture of the concepts.

Model averaging is empirically important in certain applications. For example, in statistical language modeling or speech recognition, one “smooths” a trigram model with a bigram model and smooths the bigram model with a unigram model. This smoothing is essential for minimizing the cross entropy between, say, the model and a set of corpus of newspaper sentences. It turns out that smoothing in statistical language modeling is more naturally formulated as model averaging than as model selection. A smoothed language model is very large – it contains a full trigram model, a full bigram model and a full unigram model as parts. If one uses MDL to select a single model, a language model, after a model parameters with maximum likelihood, the resulting structure is much smaller than that of a smoothed trigram model. Furthermore, the MDL model performs quite badly. However, a smoothed trigram model is theoretically derived as a compact representation of a Bayesian mixture of an exponential number of (smaller) suffix tree models [10].

Model averaging can also be applied to decision trees. A common method of constructing decision trees is to first build an overly large tree which over-fits the training data. Then one prunes the tree in such a way as to get a smaller tree that does not over-fit the data [11, 5]. An alternative to pruning is to construct a weighted mixture of the subtrees of the original over-fit tree. It is possible to construct a concise representation of a weighting over exponentially many different subtrees [3, 9, 4].

This paper proves a new PAC-Bayesian bound giving a bound on the generalization error of weighted mixtures. A weighted mixture which gives too much weight to models with low prior probability will over-fit the

Seminal paper, that contains the bound stated earlier today.

Since then, various bounds published :

- tighter,
- with less assumptions (i.i.d, bounded loss),
- easier to optimize,
- ...



Seeger, M. (2002). PAC-Bayesian generalisation error bounds for Gaussian process classification. *Journal of Machine Learning Research*.



Maurer, A. (2004). A note on the PAC-Bayesian theorem. Arxiv preprint arXiv :cs/0411099.

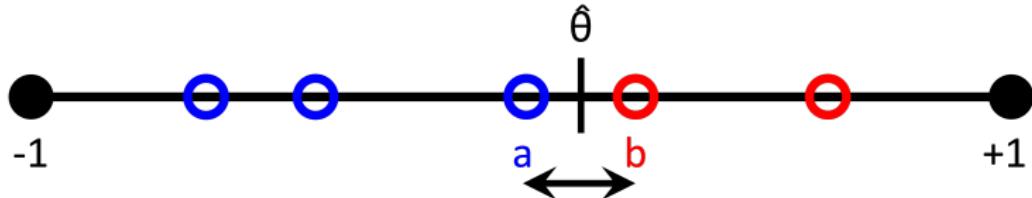


Tolstikhin, I. and Seldin, Y. (2013). PAC-Bayes-empirical-Bernstein inequality. *NeurIPS*.

Tolstikhin and Seldin's PAC-Bayes bound, 2013

With proba. at least $1 - \delta$, for any ρ ,

$$\begin{aligned} \mathbb{E}_{\theta \sim \rho}[R(\theta)] &\leq \mathbb{E}_{\theta \sim \rho}[R_n(\theta)] \\ &+ \sqrt{2\mathbb{E}_{\theta \sim \rho}[R_n(\theta)] \frac{\text{KL}(\rho\|\pi) + \log \frac{2\sqrt{n}}{\delta}}{n}} \\ &+ 2\frac{\text{KL}(\rho\|\pi) + \log \frac{2\sqrt{n}}{\delta}}{n}. \end{aligned}$$



$$\mathbb{E}_{\theta \sim \rho}[R(\theta)] \leq \sqrt{\frac{\log \frac{2}{b-a} + \log \left(\frac{2\sqrt{n}}{\delta} \right)}{2n}}.$$

$$\mathbb{E}_{\theta \sim \rho}[R(\theta)] \leq 2 \frac{\log \frac{2}{b-a} + \log \left(\frac{2\sqrt{n}}{\delta} \right)}{n}.$$

Bound on $\mathbb{E}_{\theta \sim \rho}[R(\theta)] \rightarrow$ bound on $R[\mathbb{E}_{\theta \sim \rho}(\theta)]$:



Germain, P., Lacasse, A., Laviolette, F., Marchand, M. and Roy, J.-F. (2015). Risk bounds for the majority vote : from a PAC-Bayesian analysis to a learning algorithm. *Journal of Machine Learning Research*.



Masegosa, A., Lorenzen, S., Iglesias, C. and Seldin, Y. (2020). Second order PAC-Bayesian bounds for the weighted majority vote. *NeurIPS*.

Tight bound that allows to recover all the above, and more :

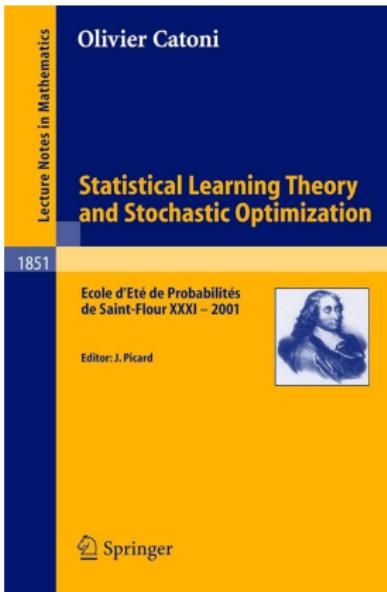


Germain, P., Lacasse, A., Laviolette, F. and Marchand, M. (2009). PAC-Bayesian learning of linear classifiers. *ICML*.

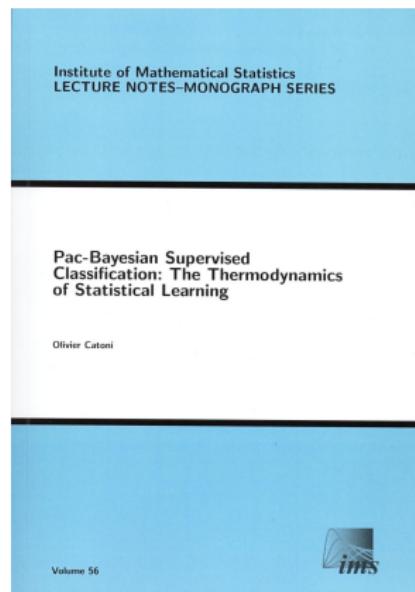


François Laviolette (1962-2021).

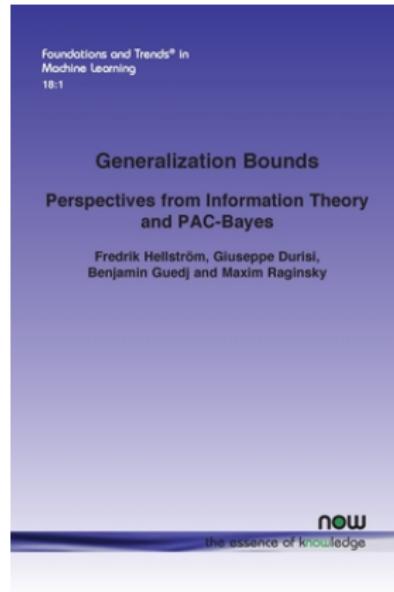
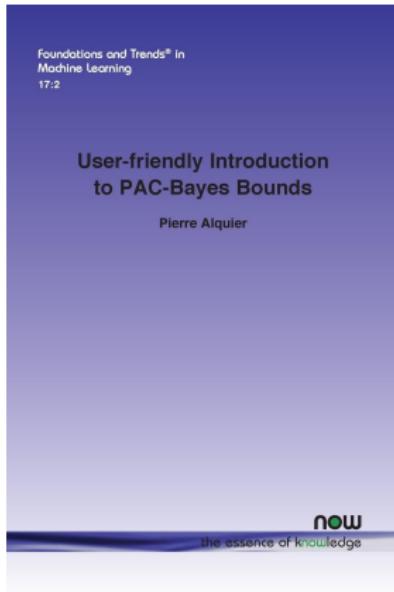
Classical references :



Connections with information theory
and MDL, oracle inequalities, rates of
convergence.



Tighter bounds, oracle inequalities,
applications to Support Vector
Machines.



Both available on arXiv...

1 PAC-Bayes bounds : introduction

- Generalization bounds and PAC-Bayes
- Minimization of the PAC-Bayes bound
- References

2 Relevance of PAC-Bayes in the post-Bayes community

- Rates of convergence
- Analysis of generalized posteriors
- Mutual Information bounds : optimizing the prior

Catoni's PAC-Bayes bound, 2003

Fix $\lambda > 0$ and π . With proba. at least $1 - \delta$ on \mathcal{S} , for any $\hat{\rho}$,

$$\mathbb{E}_{\theta \sim \hat{\rho}}[R(\theta)] \leq \mathbb{E}_{\theta \sim \hat{\rho}}[R_n(\theta)] + \frac{\text{KL}(\hat{\rho} \parallel \pi) + \log \frac{1}{\delta}}{\lambda} + \frac{\lambda}{8n}.$$

Catoni's PAC-Bayes bound in expectation, 2003

- Fix $\lambda > 0$, π and a randomized estimator $\hat{\rho}$.

$$\mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}}[R(\theta)] \right] \leq \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}}[R_n(\theta)] + \frac{\text{KL}(\hat{\rho} \parallel \pi)}{\lambda} + \frac{\lambda}{8n} \right].$$



Catoni, O. (2003). A PAC-Bayesian approach to adaptive classification. Preprint LPMA 840.



Zhang, T. (2006). Information-theoretic upper and lower bounds for statistical estimation IEEE Transactions on Information Theory.

Reminder – Gibbs posterior

$$\hat{\pi}_\lambda = \arg \min_{\rho \in \mathcal{M}(\Theta)} \left\{ \mathbb{E}_{\theta \sim \rho} [R_n(\theta)] + \frac{\text{KL}(\rho \| \pi)}{\lambda} \right\}$$

$$\hat{\pi}_\lambda(d\theta) = \frac{\exp(-\lambda R_n(\theta))}{\mathbb{E}_{\vartheta \sim \pi} [\exp(-\lambda R_n(\vartheta))]} \pi(d\theta).$$

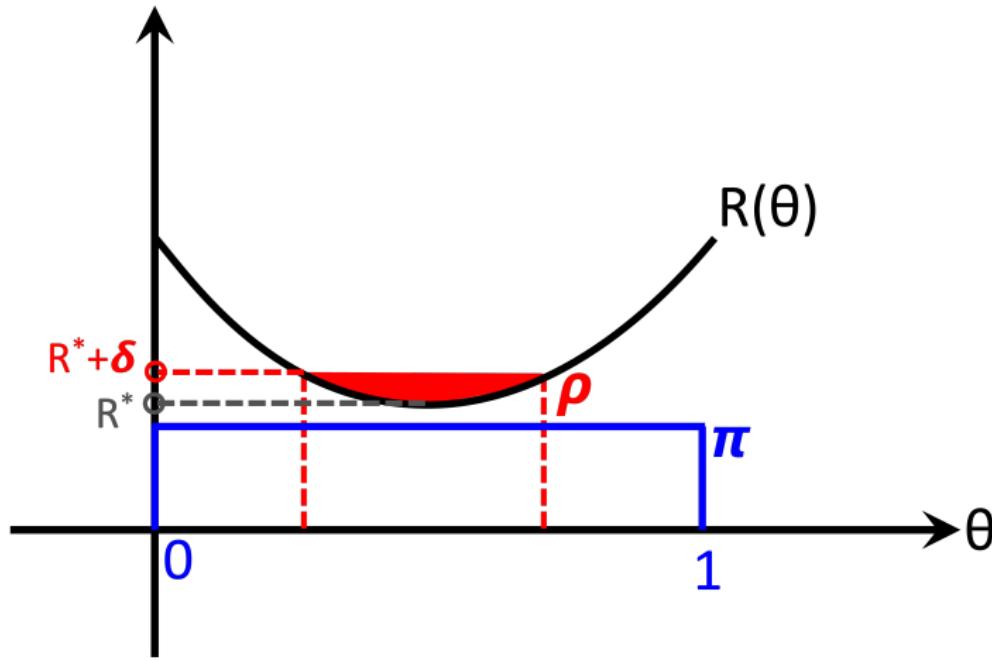
Consequence of PAC-Bayes bound in expectation :

Catoni's PAC-Bayes oracle bound, 2003

- Fix $\lambda > 0$, π , and let $\hat{\pi}_\lambda$ be the Gibbs posterior.

$$\mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\pi}_\lambda} [R(\theta)] \right] \leq \inf_{\rho \in \mathcal{M}(\Theta)} \left[\mathbb{E}_{\theta \sim \rho} [R(\theta)] + \frac{\text{KL}(\rho \| \pi)}{\lambda} + \frac{\lambda}{8n} \right].$$

$\rho = \pi_\delta := \text{restriction of } \pi \text{ to } \{\theta : R(\theta) \leq R^* + \delta\}.$



Definition : the **prior mass condition** is satisfied if there are $C, d > 0$ such that, for any $\delta > 0$ small enough,

$$\log \frac{1}{\pi\{\theta : R(\theta) \leq R^* + \delta\}} \leq d \log \frac{C}{\delta}.$$

Theorem - excess risk bound

- Assume the prior mass condition with $C, d > 0$.
- Fix $\lambda = \sqrt{n/d} \log(d/n)$, and let $\hat{\pi}_\lambda$ be the Gibbs posterior.

$$\mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\pi}_\lambda} [R(\theta)] \right] \leq R^* + \mathcal{O} \left(\sqrt{\frac{d}{n}} \log \frac{n}{d} \right).$$

Conditions for faster rates : see Catoni's book...

1 PAC-Bayes bounds : introduction

- Generalization bounds and PAC-Bayes
- Minimization of the PAC-Bayes bound
- References

2 Relevance of PAC-Bayes in the post-Bayes community

- Rates of convergence
- Analysis of generalized posteriors
- Mutual Information bounds : optimizing the prior

$$\begin{aligned}
 \hat{\theta}_{\text{MLE}} &= \arg \max_{\theta \in \Theta} \prod_{i=1}^n p_\theta(X_i) \\
 &= \arg \max_{\theta \in \Theta} \frac{\prod_{i=1}^n p_\theta(X_i)}{\prod_{i=1}^n p_{\theta_0}(X_i)} \\
 &= \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log \frac{p_{\theta_0}(X_i)}{p_\theta(X_i)}.
 \end{aligned}$$

The MLE can be seen a special case of ERM with the risk

$$R_n(\theta) := \frac{1}{n} \sum_{i=1}^n \log \frac{p_{\theta_0}(X_i)}{p_\theta(X_i)} \xrightarrow[n \rightarrow \infty]{a.s.} KL(P_{\theta_0} \| P_\theta) =: R(\theta).$$

Notation : “log-likelihood ratio”

$$LR_n(\theta_0, \theta) := \frac{1}{n} \sum_{i=1}^n \log \frac{p_{\theta_0}(X_i)}{p_\theta(X_i)}.$$

As $LR_n(\theta_0, \theta)$ is not bounded in general, we cannot apply McAllester's bound.

PAC-Bayes bound for statistical inference

Fix $\alpha \in (0, 1)$ and a prior π ,

$$\mathbb{E}_{\mathcal{S}} \mathbb{E}_{\hat{\theta}} D_{\alpha}(P_{\hat{\theta}} \| P_{\theta_0}) \leq \frac{\mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\hat{\theta} \sim \hat{\rho}} \left[\alpha LR_n(\theta_0, \hat{\theta}) \right] + \frac{KL(\hat{\rho} \| \pi)}{n} \right]}{1 - \alpha}$$

where D_{α} is the Rényi divergence.



Bhattacharya, A., Pati, D. and Yang, Y. (2019). Bayesian fractional posteriors. *The Annals of Statistics*.

$$\mathbb{E}_{\mathcal{S}} \mathbb{E}_{\hat{\theta}} D_{\alpha}(P_{\hat{\theta}} \| P_{\theta_0}) \leq \frac{\mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\hat{\theta} \sim \hat{\rho}} \left[\alpha LR_n(\theta_0, \hat{\theta}) \right] + \frac{KL(\hat{\rho} \| \pi)}{n} \right]}{1 - \alpha}.$$

The right-hand side is minimized by

$$\begin{aligned}\hat{\rho}(d\theta) &\propto \exp(-\alpha n LR_n(\theta_0, \theta)) \pi(d\theta) \\ &= \left(\prod_{i=1}^n p_{\theta}(X_i) \right)^{\alpha} \pi(d\theta).\end{aligned}$$

Used by



Alquier, P. and Ridgway, J. (2020). Concentration of tempered posteriors and of their variational approximations. *The Annals of Statistics*.

to prove rates of convergence for tempered posteriors and variational approximations.

1 PAC-Bayes bounds : introduction

- Generalization bounds and PAC-Bayes
- Minimization of the PAC-Bayes bound
- References

2 Relevance of PAC-Bayes in the post-Bayes community

- Rates of convergence
- Analysis of generalized posteriors
- Mutual Information bounds : optimizing the prior

Reminder – Catoni's PAC-Bayes bound, 2003

Fix $\lambda > 0$ and π . With proba. at least $1 - \delta$ on \mathcal{S} , for any randomized estimator $\hat{\rho}$,

$$\mathbb{E}_{\theta \sim \hat{\rho}}[R(\theta)] \leq \mathbb{E}_{\theta \sim \hat{\rho}}[R_n(\theta)] + \frac{\text{KL}(\hat{\rho} \| \pi) + \log \frac{1}{\delta}}{\lambda} + \frac{\lambda}{8n}.$$

For λ and π are fixed, this motivated the introduction of **the Gibbs posterior** $\hat{\rho} = \hat{\pi}_\lambda$, that minimizes the r.h.s. Then, we applied the bound in expectation to derive rates of convergence :

$$\mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\pi}_\lambda}[R(\theta)] \right] \leq \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\pi}_\lambda}[R_n(\theta)] + \frac{\text{KL}(\hat{\pi}_\lambda \| \pi)}{\lambda} + \frac{\lambda}{8n} \right].$$

But... why did we keep the same λ and π ?

PAC-Bayes bound in expectation – v2.0

- Fix $\Lambda > 0$, Π and the randomized estimator $\hat{\rho}$ (for example $\hat{\rho} = \hat{\pi}_\lambda$).

$$\mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R(\theta)] \right] \leq \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R_n(\theta)] + \frac{\text{KL}(\hat{\rho} \| \Pi)}{\Lambda} + \frac{\Lambda}{8n} \right].$$

Thus,

$$\begin{aligned} & \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R(\theta)] \right] \\ & \leq \inf_{\Lambda > 0} \inf_{\Pi \in \mathcal{M}(\Theta)} \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R_n(\theta)] + \frac{\text{KL}(\hat{\rho} \| \Pi)}{\Lambda} + \frac{\Lambda}{8n} \right] \\ & = \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R_n(\theta)] \right] + \inf_{\Lambda > 0} \inf_{\Pi \in \mathcal{M}(\Theta)} \mathbb{E}_{\mathcal{S}} \left[\frac{\text{KL}(\hat{\rho} \| \Pi)}{\Lambda} + \frac{\Lambda}{8n} \right]. \end{aligned}$$

$$\mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R_n(\theta)] \right] + \inf_{\Lambda > 0} \inf_{\Pi \in \mathcal{M}(\Theta)} \mathbb{E}_{\mathcal{S}} \left[\frac{\text{KL}(\hat{\rho} \| \Pi)}{\Lambda} + \frac{\Lambda}{8n} \right],$$

the infimum is reached, as shown by :



Catoni, O. (2007). *PAC-Bayesian supervised learning : the thermodynamics of statistical learning.*
IMS lecture notes – monograph series.

$$\mathbb{E}_{\mathcal{S}} \text{KL}(\hat{\rho} \| \Pi) = \mathbb{E}_{\mathcal{S}} \text{KL}(\hat{\rho} \| \mathbb{E}_{\mathcal{S}} \hat{\rho}) + \underbrace{\text{KL}(\mathbb{E}_{\mathcal{S}} \hat{\rho} \| \Pi)}_{=0 \text{ if } \Pi = \mathbb{E}_{\mathcal{S}} \hat{\rho}}.$$

- $\mathbb{E}_{\mathcal{S}} \hat{\rho} \in \mathcal{M}(\Theta)$ defined by $[\mathbb{E}_{\mathcal{S}} \hat{\rho}](E) = \mathbb{E}_{\mathcal{S}} [\hat{\rho}(E)]$.
- the first term in the r.h.s. has a nice interpretation...

Let $(U, V) \sim P$. Let P_U and P_V denote their marginals. If U and V were independent, $P = P_U \otimes P_V$.

Mutual information between two random variables

$$\mathcal{I}(U, V) := \text{KL}(P \| P_U \otimes P_V).$$

Note : $\mathcal{I}(U, V)$ depends on the distribution P of (U, V) , not on (U, V) . This is confusing... remember that $\mathbb{E}(U)$ is not a function of U !

Proposition

$$\mathcal{I}(U, V) = \mathbb{E}_U \left[\text{KL}(P_{V|U} \| P_V) \right].$$

Thus,

$$\mathbb{E}_{\mathcal{S}} \text{KL}(\hat{\rho} \| \Pi) = \underbrace{\mathbb{E}_{\mathcal{S}} \text{KL}(\hat{\rho} \| \mathbb{E}_{\mathcal{S}} \hat{\rho})}_{=: \mathcal{I}(\theta, \mathcal{S})} + \underbrace{\text{KL}(\mathbb{E}_{\mathcal{S}} \hat{\rho} \| \Pi)}_{=0 \text{ if } \Pi = \mathbb{E}_{\mathcal{S}} \hat{\rho}}.$$

$$\begin{aligned}
 & \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R(\theta)] \right] \\
 & \leq \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R_n(\theta)] \right] + \inf_{\Lambda > 0} \inf_{\Pi \in \mathcal{M}(\Theta)} \mathbb{E}_{\mathcal{S}} \left[\frac{\text{KL}(\hat{\rho} \parallel \Pi)}{\Lambda} + \frac{\Lambda}{8n} \right] \\
 & = \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R_n(\theta)] \right] + \inf_{\Lambda > 0} \left[\frac{\mathcal{I}(\theta, \mathcal{S})}{\Lambda} + \frac{\Lambda}{8n} \right].
 \end{aligned}$$

Mutual information bound

$$\mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R(\theta)] \right] \leq \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R_n(\theta)] \right] + \sqrt{\frac{\mathcal{I}(\theta, \mathcal{S})}{2n}}.$$

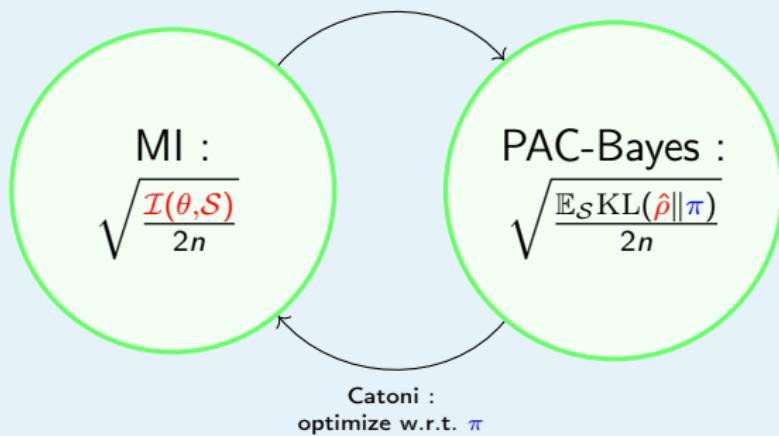


Russo, D. and Zou, J. (2019). How much does your data exploration overfit? controlling bias via information usage. *IEEE Transactions on Information Theory*.

Mutual information bound

$$\mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R(\theta)] \right] \leq \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\theta \sim \hat{\rho}} [R_n(\theta)] \right] + \sqrt{\frac{\mathcal{I}(\theta, \mathcal{S})}{2n}}.$$

$$\mathcal{I}(\theta, \mathcal{S}) = \mathbb{E}_{\mathcal{S}} \text{KL}(\hat{\rho} \| \mathbb{E}_{\mathcal{S}} \hat{\rho}) \leq \mathbb{E}_{\mathcal{S}} \text{KL}(\hat{\rho} \| \pi)$$



Example

Assume a prior mass condition. Apply ~~PAC~~-Bayes mutual information bound :

$$\mathbb{E}_S \left[\mathbb{E}_{\theta \sim \hat{\pi}_\lambda} [R(\theta)] \right] \leq R^* + O \left(\sqrt{\frac{d}{n} \log \frac{n}{d}} \right).$$

For classification : see Catoni's book.

For log-likelihood and tempered posteriors :



EL Mahdi Khribch, Pierre Alquier (2024). Convergence of Statistical Estimators via Mutual Information Bounds. *Preprint arXiv :2412.18539*.

THANK YOU !

Contact information :

- contact : alquier@essec.edu
- webpage : <https://pierrealquier.github.io/>

Many thanks to Richard III Cariño who helped with the drawings !