

PAC-Bayesian Theory Meets Bayesian Inference

Authored by:

Simon Lacoste-Julien
Pascal Germain
Francis Bach
Alexandre Lacoste

Abstract

We exhibit a strong link between frequentist PAC-Bayesian bounds and the Bayesian marginal likelihood. That is, for the negative log-likelihood loss function, we show that the minimization of PAC-Bayesian generalization bounds maximizes the Bayesian marginal likelihood. This provides an alternative explanation to the Bayesian Occam's razor criteria, under the assumption that the data is generated by an i.i.d. distribution. Moreover, as the negative log-likelihood is an unbounded loss function, we motivate and propose a PAC-Bayesian theorem tailored for the sub-gamma loss family, and we show that our approach is sound on classical Bayesian linear regression tasks.

1 Paper Body

Since its early beginning [24, 34], the PAC-Bayesian theory claims to provide PAC guarantees to Bayesian algorithms? (McAllester [24]). However, despite the amount of work dedicated to this statistical learning theory?many authors improved the initial results [8, 21, 25, 30, 35] and/or generalized them for various machine learning setups [4, 12, 15, 20, 28, 31, 32, 33]?it is mostly used as a frequentist method. That is, under the assumptions that the learning samples are i.i.d.-generated by a data-distribution, this theory expresses probably approximately correct (PAC) bounds on the generalization risk. In other words, with probability 1, the generalization risk is at most " away from the training risk. The Bayesian side of PAC-Bayes comes mostly from the fact that these bounds are expressed on the averaging/aggregation/ensemble of multiple predictors (weighted by a posterior distribution) and incorporate prior knowledge. Although it is still sometimes referred as a theory that bridges the Bayesian and frequentist approach [e.g., 16], it has been merely used to justify Bayesian methods until now.¹ In this work, we provide a direct connection between Bayesian inference techniques [summarized by 5, 13] and PAC-Bayesian risk bounds in a general setup. Our study is based on a simple but insightful connection between

the Bayesian marginal likelihood and PAC-Bayesian bounds (previously mentioned by Grünwald [14]) obtained by considering the negative log-likelihood loss function (Section 3). By doing so, we provide an alternative explanation for the Bayesian Occam’s razor criteria [18, 22] in the context of model selection, expressed as the complexity-accuracy trade-off appearing in most PAC-Bayesian results. In Section 4, we extend PAC-Bayes theorems to regression problems with unbounded loss, adapted to the negative log-likelihood loss function. Finally, we study the Bayesian model selection from a PAC-Bayesian perspective (Section 5), and illustrate our finding on classical Bayesian regression tasks (Section 6).

2

PAC-Bayesian Theory

We denote the learning sample $(X, Y) = \{(x_i, y_i)\}_{i=1}^n \subset (X \times Y)^n$, that contains n input-output pairs. The main assumption of frequentist learning theories including PAC-Bayes is that (X, Y) is 1

Some existing connections [3, 6, 14, 19, 29, 30, 36] are discussed in Appendix A.1.

30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

randomly sampled from a data generating distribution that we denote D . Thus, we denote $(X, Y) \in D^n$ the i.i.d. observation of n elements. From a frequentist perspective, we consider in this work loss functions $\ell : F \times X \times Y \rightarrow \mathbb{R}$, where F is a (discrete or continuous) set of predictors $f : X \rightarrow Y$, and we write the empirical risk on the sample (X, Y) and the generalization error on distribution D as $L_{X,Y}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f, x_i, y_i)$; $L_D(f) = \mathbb{E} \ell(f, x, y)$. The PAC-Bayesian theory [24, 25] studies an averaging of the above losses according to a posterior distribution π over F . That is, it provides probably approximately correct generalization bounds on the (unknown) quantity $L_D(\pi) = \mathbb{E}_{\pi} L_D(f) = \mathbb{E}_{\pi} \mathbb{E}_{(x,y) \sim D} \ell(f, x, y)$, given the empirical estimate $L_{X,Y}(f)$ and some other parameters. Among these, most PAC-Bayesian theorems rely on the Kullback-Leibler divergence $KL(\pi \| \eta) = \mathbb{E}_{\pi} \ln[\pi(f)/\eta(f)]$ between a prior distribution η over F specified before seeing the learning sample X, Y and the posterior π typically obtained by feeding a learning process with (X, Y) .

Two appealing aspects of PAC-Bayesian theorems are that they provide data-driven generalization bounds that are computed on the training sample (i.e., they do not rely on a testing sample), and that they are uniformly valid for all π over F . This explains why many works study them as model selection criteria or as an inspiration for learning algorithm conception. Theorem 1, due to Catoni [8], has been used to derive or study learning algorithms [10, 17, 26, 27]. Theorem 1 (Catoni [8]). Given a distribution D over $X \times Y$, a hypothesis set F , a loss function $\ell : F \times X \times Y \rightarrow [0, 1]$, a prior distribution η over F , a real number $\beta \in (0, 1]$, and a real number $\gamma \geq 0$, with probability at least $1 - \beta$ over the choice of $(X, Y) \in D^n$, we have $\frac{1}{n} L_{X,Y}(\pi) \leq \frac{1}{n} L_{X,Y}(\eta) + \beta \ln \frac{1}{\beta} + \gamma$ on $F : \mathbb{E} L_D(\pi) \leq \frac{1}{n} L_{X,Y}(\eta) + \beta \ln \frac{1}{\beta} + \gamma$. Theorem 1 is limited to loss functions mapping to the range $[0, 1]$. Through a straightforward rescaling we

can extend it to any bounded loss, i.e., $\ell : \mathcal{F} \times \mathcal{X} \times \mathcal{Y} \rightarrow [a, b]$, where $[a, b] \subset \mathbb{R}$. This is done by using $\ell := (b - a)\ell_0$ and with the rescaled loss function $\ell_0(f, x, y) := (\ell(f, x, y) - a)/(b - a) \in [0, 1]$. After few arithmetic manipulations, we can rewrite Equation (1) as $\mathbb{E} \ell(f, X, Y) = \mathbb{E} \ell_0(f, X, Y) + a$ and $\mathbb{E} \ell_0(f, X, Y) = \int \int \ell_0(f, x, y) p(x, y) dx dy$.

$$\frac{1}{n} \sum_{i=1}^n \ell(f, x_i, y_i)$$

From an algorithm design perspective, Equation (2) suggests optimizing a trade-off between the empirical expected loss and the Kullback-Leibler divergence. Indeed, for fixed λ, X, Y, n , and \mathcal{F} , minimizing Equation (2) is equivalent to find the distribution q that minimizes $\mathbb{E} \ell(f, X, Y) + \lambda \text{KL}(q \| p)$. (3)

It is well known [1, 8, 10, 21] that the optimal Gibbs posterior q^* is given by $q^*(f) = \frac{1}{Z(\lambda, X, Y)} \int \exp(-\lambda \ell(f, x, y)) p(x, y) dx dy$,

$$\begin{aligned} & \int \exp(-\lambda \ell(f, x, y)) p(x, y) dx dy \\ & \int \exp(-\lambda \ell(f, x, y)) p(x, y) dx dy \end{aligned} \quad (4)$$

where $Z(\lambda, X, Y)$ is a normalization term. Notice that the constant of Equation (1) is now absorbed in the loss function as the rescaling factor setting the trade-off between the expected empirical loss and $\text{KL}(q \| p)$.

3

Bridging Bayes and PAC-Bayes

In this section, we show that by choosing the negative log-likelihood loss function, minimizing the PAC-Bayes bound is equivalent to maximizing the Bayesian marginal likelihood. To obtain this result, we first consider the Bayesian approach that starts by defining a prior $p(\theta)$ over the set of possible model parameters θ . This induces a set of probabilistic estimators $f_\theta \in \mathcal{F}$, mapping x to a probability distribution over \mathcal{Y} . Then, we can estimate the likelihood of observing y given x and θ , i.e., $p(y|x, \theta) = \int p(y|x, \theta) p(\theta) d\theta$. Using Bayes' rule, we obtain the posterior $p(\theta|X, Y) = \frac{p(\theta) \prod_{i=1}^n p(y_i|x_i, \theta)}{\int p(\theta) \prod_{i=1}^n p(y_i|x_i, \theta) d\theta}$. (5) $p(Y|X) = \int p(Y|X, \theta) p(\theta) d\theta$ where $p(Y|X, \theta) = \prod_{i=1}^n p(y_i|x_i, \theta)$ and $p(Y|X) = \int p(Y|X, \theta) p(\theta) d\theta$. 2

To stay aligned with the PAC-Bayesian setup, we only consider the discriminative case in this paper. One can extend to the generative setup by considering the likelihood of the form $p(y, x|\theta)$ instead.

2

To bridge the Bayesian approach with the PAC-Bayesian framework, we consider the negative log-likelihood loss function [3], denoted ℓ_{nll} and defined by $\ell_{\text{nll}}(f, x, y) = -\ln p(y|x, \theta)$.

$$\ell_{\text{nll}}(f, x, y) = -\ln p(y|x, \theta)$$

Then, we can relate the empirical loss $L_{X,Y}$ of a predictor to its likelihood:

n

$$\frac{1}{n} \sum_{i=1}^n \ell_{X,Y}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\theta, x_i, y_i) \quad \text{or, the other way around,}$$

$$\frac{1}{n} \sum_{i=1}^n \ln p(y_i - x_i, \theta) = \frac{1}{n} \sum_{i=1}^n \ln p(Y - X, \theta) = \mathbb{E} \ell_{X,Y}(\theta)$$

(7)

Unfortunately, existing PAC-Bayesian theorems work with bounded loss functions or in very specific contexts [e.g., 9, 36], and $\ell_{X,Y}$ spans the whole real axis in its general form. In Section 4, we explore PAC-Bayes bounds for unbounded losses. Meanwhile, we consider priors with bounded likelihood. This can be done by assigning a prior of zero to any θ yielding $\ln p(y - x, \theta) \leq 0$ for $\theta \in [a, b]$. Now, using Equation (7) in the optimal posterior (Equation 4) simplifies to

$$\frac{1}{Z} \int_{\Theta} \exp\left(\frac{1}{n} \sum_{i=1}^n \ln p(y_i - x_i, \theta)\right) p(\theta) d\theta = \frac{1}{Z} \int_{\Theta} \exp\left(\frac{1}{n} \sum_{i=1}^n \ln p(y_i - x_i, \theta)\right) d\theta$$

where the normalization constant Z corresponds to the Bayesian marginal likelihood: $Z = \int_{\Theta} \exp\left(\frac{1}{n} \sum_{i=1}^n \ln p(y_i - x_i, \theta)\right) d\theta$.

(8)

(9)

?

This shows that the optimal PAC-Bayes posterior given by the generalization bound of Theorem 1 coincides with the Bayesian posterior, when one chooses $\ell_{X,Y}$ as loss function and $\theta \in [a, b]$ (as in Equation 2). Moreover, using the posterior of Equation (8) inside Equation (3), we obtain $\mathbb{E} \ell_{X,Y}(\theta) + \text{KL}(\theta \parallel p) \leq \frac{1}{n} \sum_{i=1}^n \ln p(y_i - x_i, \theta) + \text{KL}(\theta \parallel p) \leq \frac{1}{n} \sum_{i=1}^n \ln p(y_i - x_i, \theta) + \ln Z$.

(10)

?

In other words, minimizing the PAC-Bayes bound is equivalent to maximizing the marginal likelihood. Thus, from the PAC-Bayesian standpoint, the latter encodes a trade-off between the averaged negative log-likelihood loss function and the prior-posterior Kullback-Leibler divergence. Note that Equation (10) has been mentioned by Grunwald [14], based on an earlier observation of Zhang [36]. However, the PAC-Bayesian theorems proposed by the latter do not bound the generalization loss directly, as the ‘classical’ PAC-Bayesian results [8, 24, 29] that we extend to regression in forthcoming Section 4 (see the corresponding remarks in Appendix A.1). We conclude this section by proposing a compact form of Theorem 1 by expressing it in terms of the marginal likelihood, as a direct consequence of Equation (10). Corollary 2. Given a data distribution D , a parameter set Θ , a prior distribution p over Θ , a $\theta \in [a, b]$, if $\ell_{X,Y}$ lies in $[a, b]$,

we have, with probability at least $1 - \delta$ over the choice of $(X, Y) \sim D^n$, $\mathbb{E} L(\hat{f}) \leq \inf_{f \in \mathcal{F}} \mathbb{E} L(f) + \sqrt{\frac{2 \log(2/\delta)}{n}}$

where \mathbb{E} is the Gibbs optimal posterior (Eq. 8) and $L(f)$ is the marginal likelihood (Eq. 9).

In Section 5, we exploit the link between PAC-Bayesian bounds and Bayesian marginal likelihood to expose similarities between both frameworks in the context of model selection. Beforehand, next Section 4 extends the PAC-Bayesian generalization guarantees to unbounded loss functions. This is mandatory to make our study fully valid, as the negative log-likelihood loss function is in general unbounded (as well as other common regression losses).

4

PAC-Bayesian Bounds for Regression

This section aims to extend the PAC-Bayesian results of Section 3 to real valued unbounded loss. These results are used in forthcoming sections to study ℓ_2 , but they are valid for broader classes of loss functions. Importantly, our new results are focused on regression problems, as opposed to the usual PAC-Bayesian classification framework. The new bounds are obtained through a recent theorem of Alquier et al. [1], stated below (we provide a proof in Appendix A.2 for completeness). Theorem 3 (Alquier et al. [1]). Given a distribution D over $X \times Y$, a hypothesis set \mathcal{F} , a loss function $\ell : \mathcal{F} \times X \times Y \rightarrow \mathbb{R}$, a prior distribution π over \mathcal{F} , $a \in (0, 1]$, and a real number $\delta \in (0, 1)$, with probability at least $1 - \delta$ over the choice of $(X, Y) \sim D^n$, we have $\mathbb{E} L(\hat{f}) \leq \inf_{f \in \mathcal{F}} \mathbb{E} L(f) + \frac{1}{a} \sqrt{\frac{2 \log(2/\delta)}{n}} + \ln \frac{1}{\pi(\mathcal{F})} + \frac{1}{a} \sqrt{\frac{2 \log(2/\delta)}{n}}$, (11)

where

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

Alquier et al. used Theorem 3 to design a learning algorithm for $\{0, 1\}$ -valued classification losses. Indeed, a bounded loss function $\ell : \mathcal{F} \times X \times Y \rightarrow [a, b]$ can be used along with Theorem 3 by applying the Hoeffding's lemma to Equation (12), that gives $\mathbb{E} L(f) \leq \frac{(b-a)^2}{2n} + \frac{1}{a} \sqrt{\frac{2 \log(2/\delta)}{n}}$. More specifically, with $a = 0$, we obtain the following bound $\mathbb{E} L(\hat{f}) \leq \inf_{f \in \mathcal{F}} \mathbb{E} L(f) + \frac{1}{a} \sqrt{\frac{2 \log(2/\delta)}{n}} + \ln \frac{1}{\pi(\mathcal{F})} + \frac{1}{a} \sqrt{\frac{2 \log(2/\delta)}{n}}$. (13)

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

$\mathbb{E} L(f) = \mathbb{E} \sum_{i=1}^n \ell(f, X_i, Y_i)$

Under this sub-gamma assumption, we obtain the following new result, which is necessary to study linear regression in the next sections. 4

Corollary 5. Given D, F, ℓ, γ and defined in the statement of Theorem 3, if the loss is sub-gamma with variance factor s^2 and scale $c \leq 1$, we have, with probability at least $1 - \delta$ over $(X, Y) \sim D^n$, $\ell(F : E \text{LD}(\ell)) \leq E \text{Lb}_{X,Y}(\ell) + n \text{KL}(\ell \| \ell^*) + \ln \frac{1}{1 - \delta} + 2(11c) s^2 \cdot \ell^*(F)$

As a special case, with $\ell := \ell_{\text{nl}}$ and $\ell^* := \ell^*_{\text{nl}}$ (Equation 8), we have $E \text{L}(\ell_{\text{nl}}) \leq$

$$\begin{aligned} & \ell^*_{\text{nl}}(F) \\ & s^2 \frac{2}{1 - c} \\ & \frac{1}{n} \\ & (18) \\ & \ln \text{Lb}_{X,Y}(\ell) \end{aligned}$$

Proof. Following the same path as in the proof of Corollary 4 (with P_n, Q_n $P_n(\cdot) = \int \delta_{\ell} dP_n = \int \delta_{\ell} dP_n = \int \delta_{\ell} dP_n$), we have $\ell^*(F) =$

$$\begin{aligned} & \frac{2}{n} \text{Lb}_{X,Y}(\ell) \\ & c) \\ & = \\ & n s^2 \frac{2}{1 - c} \end{aligned}$$

where the inequality comes from the sub-gamma loss assumption, with $\ell^*(F) \leq (0, 1/c)$. Squared loss. The parameters s and c of Corollary 5 rely on the chosen loss function and prior, and the assumptions concerning the data distribution. As an example, consider a regression problem where $X \sim \mathbb{R}^d$, a family of linear predictors $f_w(x) = w \cdot x$, with $w \in \mathbb{R}^d$, and a Gaussian prior $N(0, \frac{1}{2}I)$. Let us assume that the input examples are generated by $x \sim N(0, \frac{1}{2}I)$ with label $y = w \cdot x + \epsilon$, where $w \in \mathbb{R}^d$ and $\epsilon \sim N(0, \frac{1}{2})$ is a Gaussian noise. Under the squared loss function $\ell(y) = y^2$,

(19) we show in Appendix A.4 that Corollary 5 is valid with $s^2 = 2$ ($\frac{1}{2}d + kw^2$) and $c = \frac{1}{2}$. As expected, the bound degrades when the noise increases $\ell(w, x, y) = (w \cdot x - y)^2$

Regression versus classification. The classical PAC-Bayesian theorems are stated in a classification context and bound the generalization error/loss of the stochastic Gibbs predictor G . In order to predict the label of an example $x \in X$, the Gibbs predictor first draws a hypothesis $h \in F$ according to G , and then returns $h(x)$. Maurer [23] shows that we can generalize PAC-Bayesian bounds on the generalization risk of the Gibbs classifier to any loss function with output between zero and one. Provided that $y \in \{0, 1\}$ and $h(x) \in [0, 1]$, a common choice is to use the 1 linear loss function $\ell(h, x, y) = |y - h(x)|$. The Gibbs generalization loss is then given by $\mathbb{E} \text{Lb}_{X,Y}(\ell) = \mathbb{E} \text{Lb}_{X,Y}(\ell)$. Many PAC-Bayesian works use $\text{RD}(G)$ as a surrogate loss to study the

zero-one classification loss of the majority vote classifier $RD(B^{??}) = \mathbb{E} \sum_{(x,y) \in D} \mathbb{I}(h(x) \neq y)$ (20) $\mathbb{E} \sum_{(x,y) \in D} \mathbb{I}(h(x) \neq y)$

where $\mathbb{I}[\cdot]$ being the indicator function. Given a distribution \mathbb{P} , an upper bound on the Gibbs risk is converted to an upper bound on the majority vote risk by $RD(B^{??}) \leq 2RD(G^{??})$ [20]. In some situations, this factor of two may be reached, i.e., $RD(B^{??}) = 2RD(G^{??})$. In other situations, we may have $RD(B^{??}) = 0$ even if $RD(G^{??}) = 1/2$ (see Germain et al. [11] for an extensive study). Indeed, these bounds obtained via the Gibbs risk are exposed to be loose and/or unrepresentative of the majority vote generalization error.³ In the current work, we study regression losses instead of classification ones. That is, the provided results express upper bounds on $\mathbb{E} \sum_{(x,y) \in D} (f(x) - y)^2$ for any (bounded, sub-Gaussian, or sub-gamma) losses. Of course, one may want to bound the regression loss of the averaged regressor $F^{??}(x) = \mathbb{E} f^{??}(x)$. In this case, if the loss function ℓ is convex (as the squared loss), Jensen's inequality gives $\mathbb{E} \sum_{(x,y) \in D} \ell(F^{??}(x) - y) \leq \mathbb{E} \sum_{(x,y) \in D} \ell(f^{??}(x) - y)$. Note that a strict inequality replaces the factor two mentioned above for the classification case, due to the non-convex indicator function of Equation (20). Now that we have generalization bounds for real-valued loss functions, we can continue our study linking PAC-Bayesian results to Bayesian inference. In the next section, we focus on model selection.

5
5

Analysis of Model Selection

We consider L distinct models $\{M_i\}_{i=1}^L$, each one defined by a set of parameters θ_i . The PACBayesian theorems naturally suggest selecting the model that is best adapted for the given task by evaluating the bound for each model $\{M_i\}_{i=1}^L$ and selecting the one with the lowest bound [2, 25, 36]. This is closely linked with the Bayesian model selection procedure, as we showed in Section 3 that minimizing the PAC-Bayes bound amounts to maximizing the marginal likelihood. Indeed, given a collection of L optimal Gibbs posteriors⁴ one for each model⁵ given by Equation (8), $p(\theta_i | X, Y, M_i) \propto \prod_{(x,y) \in D} p(y | x, \theta_i)$

$\propto \prod_{(x,y) \in D} p(y | x, \theta_i)$
for $i = 1, \dots, L$,

the Bayesian Occam's razor criteria [18, 22] chooses the one with the higher model evidence $Z \propto \int p(Y | X, M_i) p(\theta_i) d\theta_i$.

(21) (22)
 θ_i

Corollary 6 below formally links the PAC-Bayesian and the Bayesian model selection. To obtain this result, we simply use the bound of Corollary 5 L times,

together with ‘nll and Equation (10). From the union bound (a.k.a. Bonferroni inequality), it is mandatory to compute each bound with a confidence parameter of α/L , to ensure that the final conclusion is valid with probability at least $1 - \alpha$. Corollary 6. Given a data distribution D , a family of model parameters $\{\theta_i\}_{i=1}^L$ and associated priors $\{p_i\}_{i=1}^L$ where θ_i is defined over Θ , $\alpha \in (0, 1]$, if the loss is sub-gamma with parameters σ_i $i=1$ to s_2 and $c \leq 1$, then, with probability at least $1 - \alpha$ over $(X, Y) \sim D^n$, $\exists i \in \{1, \dots, L\} : E \theta_i^L \leq \text{nll}(\theta_i) + \frac{2(1+c)s_2}{n} \ln \sum_{i=1}^L Z_{X,Y,i}$.

where θ_i is the Gibbs optimal posterior (Eq. 21) and $Z_{X,Y,i}$ is the marginal likelihood (Eq. 22). Hence, under the uniform prior over the L models, choosing the one with the best model evidence is equivalent to choosing the one with the lowest PAC-Bayesian bound. Hierarchical Bayes. To perform proper inference on hyperparameters, we have to rely on the Hierarchical Bayes approach. This is done by considering an hyperprior $p(\theta)$ over the set of hyperparameters Θ . Then, the prior $p(\theta_i)$ can be conditioned on a choice of hyperparameter θ . The $p(Y - X, \theta)$ Bayes rule of Equation (5) becomes $p(\theta, \theta - X, Y) = p(\theta) p(\theta - X) \cdot p(Y - X)$.

Under the negative log-likelihood loss function, we can rewrite the results of Corollary 5 as a generalization bound on $E \theta_i^L \leq \text{nll}(\theta_i)$, where θ_i is the hyperposterior on Θ and θ the hyperprior. Indeed, Equation (18) becomes $E \theta_i^L \leq \text{nll}(\theta_i) = E \theta_i^L \leq \text{nll}(\theta_i) + \frac{2(1+c)s_2}{n} \ln E \sum_{i=1}^L Z_{X,Y,i}$. (23)

where

θ_i is the hyperposterior on Θ and θ the hyperprior.

To relate to the bound obtained in Corollary 6, we consider the case of a discrete hyperparameter set $\Theta = \{\theta_i\}_{i=1}^L$, with a uniform prior $p(\theta_i) = 1/L$ (from now on, we regard each hyperparameter θ_i as the specification of a model θ_i). Then, Equation (23) becomes $E \theta_i^L \leq \text{nll}(\theta_i) = E \theta_i^L \leq \text{nll}(\theta_i) + \frac{2(1+c)s_2}{n} \ln \sum_{i=1}^L Z_{X,Y,i}$.

where

This bound is now a function of $\sum_{i=1}^L Z_{X,Y,i}$ instead of $\max_i Z_{X,Y,i}$ as in the bound given by the ‘best’ model in Corollary 6. This yields a tighter bound, corroborating the Bayesian wisdom that model averaging performs best. Conversely, when selecting a single hyperparameter $\theta \in \Theta$, the hierarchical representation is equivalent to choosing a deterministic hyperposterior, satisfying $p(\theta_i) = 1$ and 0 for every other values. We then have $\text{KL}(\theta_i - \theta) = \text{KL}(\theta_i - \theta) + E \text{KL}(\theta_i - \theta) = \ln(L) + \text{KL}(\theta_i - \theta)$.

With the optimal posterior for the selected θ , we have $\text{nll}(\theta) = E \text{nll}(\theta) + \text{KL}(\theta - \theta) = E \text{nll}(\theta) + \text{KL}(\theta - \theta) + \ln(L)$.

Inserting this result into Equation (17), we fall back on the bound obtained in Corollary 6. Hence, by comparing the values of the bounds, one can get an estimate on the consequence of performing model selection instead of model averaging.

Linear Regression

In this section, we perform Bayesian linear regression using the parameterization of Bishop [5]. The output space is $Y := \mathbb{R}$ and, for an arbitrary input space X , we use a mapping function $\phi : X \rightarrow \mathbb{R}^d$. The model. Given $(x, y) \in X \times Y$ and model parameters $\theta := \{w, \sigma^2\}$, $w \in \mathbb{R}^d$, $\sigma^2 \in \mathbb{R}^+$, we consider the likelihood $p(y|x, w, \sigma^2) = \mathcal{N}(y|w^T \phi(x), \sigma^2)$. Thus, the negative log-likelihood loss is

$$\begin{aligned} \text{‘nll}(w, \sigma^2, x, y) &= \\ &= -\ln p(y|x, w, \sigma^2) = \\ &= -\ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y - w^T \phi(x))^2\right) \right) \\ &= \frac{1}{2\sigma^2} (y - w^T \phi(x))^2 + \frac{1}{2} \ln(2\pi\sigma^2). \end{aligned} \quad (24)$$

For a fixed σ^2 , minimizing Equation (24) is equivalent to minimizing the squared loss function of Equation (19). We also consider an isotropic Gaussian prior of mean 0 and variance σ_0^2 : $p(w) = \mathcal{N}(w|0, \sigma_0^2 I)$. For the sake of simplicity, we consider fixed parameters σ^2 and σ_0^2 . The Gibbs optimal posterior (see Equation 8) is then given by

$$\begin{aligned} p(w|X, Y, \sigma^2, \sigma_0^2) &= \\ &= \frac{p(w) \prod_{i=1}^n p(y_i|x_i, w, \sigma^2)}{\int dw p(w) \prod_{i=1}^n p(y_i|x_i, w, \sigma^2)} \\ &= \mathcal{N}(w|w_0, \Sigma_0^{-1}) \end{aligned} \quad (25)$$

where $\Sigma_0 := \sigma_0^2 I$; $w_0 := \frac{1}{n} \sum_{i=1}^n y_i \phi(x_i)$; $\Sigma_0^{-1} := \frac{1}{\sigma_0^2} I + \frac{1}{\sigma^2} \sum_{i=1}^n \phi(x_i) \phi(x_i)^T$; $w := [w_1, \dots, w_n]^T$ is the labels-vector; and the negative log marginal likelihood is $-\ln p(Y|X, \sigma^2, \sigma_0^2) =$

$$\begin{aligned} &= -\ln \int dw p(w) \prod_{i=1}^n p(y_i|x_i, w, \sigma^2) \\ &= -\ln \int dw \mathcal{N}(w|w_0, \Sigma_0^{-1}) \prod_{i=1}^n \mathcal{N}(y_i|x_i, w, \sigma^2) \\ &= -\ln \int dw \frac{1}{(2\pi)^{d/2} |\Sigma_0^{-1}|^{1/2}} \exp\left(-\frac{1}{2} w^T \Sigma_0^{-1} w + w^T \Sigma_0^{-1} w_0\right) \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y_i - w^T \phi(x_i))^2\right) \\ &= -\ln \int dw \exp\left(-\frac{1}{2} w^T \Sigma_0^{-1} w + w^T \Sigma_0^{-1} w_0 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T \phi(x_i))^2\right) \\ &= -\ln \int dw \exp\left(-\frac{1}{2} w^T \left(\Sigma_0^{-1} + \frac{1}{\sigma^2} \sum_{i=1}^n \phi(x_i) \phi(x_i)^T\right) w + w^T \left(\Sigma_0^{-1} w_0 + \frac{1}{\sigma^2} \sum_{i=1}^n y_i \phi(x_i)\right)\right) \\ &= -\ln \int dw \exp\left(-\frac{1}{2} w^T \Sigma^{-1} w + w^T \Sigma^{-1} w_0\right) \end{aligned}$$

is designed to reproduce Bishop [5, Figure 3.14], where it is explained that the marginal likelihood correctly indicates that the polynomial model of degree $d = 3$ is ‘the simplest model which gives a good explanation for the observed data’. We show that this claim is well quantified by the trade-off intrinsic to our PAC-Bayesian approach: the complexity KL term keeps increasing with the parameter $d \in \{1, 2, \dots, 7\}$, while the empirical risk drastically decreases from $d = 2$ to $d = 3$, and only slightly afterward. Moreover, we show that the generalization risk (computed on a test sample of size 1000) tends to increase with complex models (for $d \geq 4$). Empirical comparison of bound values. Figure 1c compares the values of the PAC-Bayesian bounds presented in this paper on a synthetic dataset, where each input $x \in \mathbb{R}^{20}$ is generated by a Gaussian $x \sim \mathcal{N}(0, I)$. The associated output $y \in \mathbb{R}$ is given by $y = w^\top x + \epsilon$, with $k = 12$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$, and $\sigma^2 = 19$. We perform Bayesian linear regression in the input space, i.e., $\phi(x) = x$, fixing $\sigma^2 = 100$ and $\lambda = 2$. That is, we compute the posterior of Equation (25) for training samples of sizes from 10 to 106. For each learned model, we compute the empirical negative log-likelihood loss ℓ of Equation (24), and the three PAC-Bayes bounds, with confidence parameter of $\delta = 20$. Note that this loss function is an affine transformation of the squared loss studied in Section 4 (Equation 19), i.e., $\ell(w, \phi; x, y) = 12 \ln(2\sigma^2) + 2 \|\phi(x) - y\|^2 / \sigma^2$, out that ℓ is sub-gamma with parameters $\beta(x, y)$. It turns out that ℓ is sub-gamma with parameters $\beta(x, y)$, as shown in Appendix A.6. The bounds $\mathbb{E}[\ell] + \delta \sqrt{\beta(x, y)}$ and c of Corollary 5 are computed using the above mentioned values of k, d, σ^2, x, y , leading to

1.5
model model model model
1.0 0.5
d=1 d=2 d=3 d=4
60
model d=5 model d=6 model d=7 $\sin(x)$
50 40
0.0
$\ln \mathbb{E}[\ell]$
30
$\text{KL}(\mathbb{P} \ \mathbb{Q})$ ℓ $\mathbb{E}[\ell]$ $\mathbb{E}[\ell] + \delta \sqrt{\beta(x, y)}$ c
0.5
20
1.0
$\mathbb{E}[\ell] + \delta \sqrt{\beta(x, y)}$
10
1.5 2.0 0
1 2
3 2
?
0 1
2
2

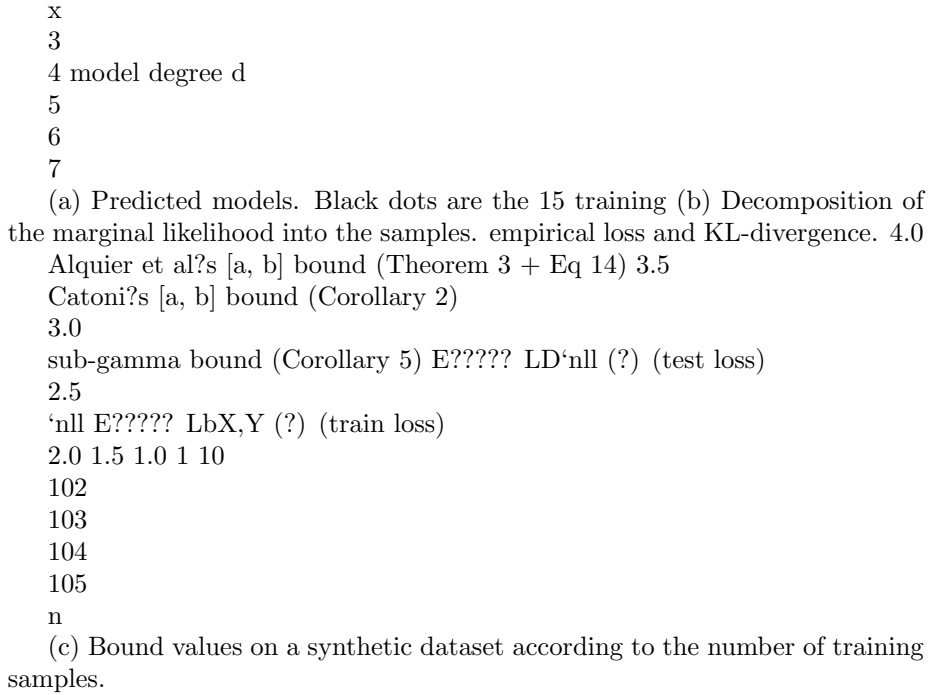


Figure 1: Model selection experiment (a-b); and comparison of bounds values (c). to $s_2 \approx 0.280$ and $c \approx 0.005$. As the two other bounds of Figure 1c are not suited for unbounded loss, we compute their value using a cropped loss $[a, b] = [1, 4]$. Different parameter values could have been chosen, sometimes leading to another picture: a large value of s degrades our sub-gamma bound, as a larger $[a, b]$ interval does for the other bounds. In the studied setting, the bound of Corollary 5 that we have developed for (unbounded) subgamma losses gives tighter guarantees than the two results for $[a, b]$ -bounded losses (up to $n=106$). However, our new bound always maintains a gap of $2(11c)s_2$ between its value and the generalization loss. The result of Corollary 2 (adapted from Catoni [8]) for bounded losses suffers from a similar gap, while having higher values than our sub-gamma result. Finally, the result of Theorem 3 (Alquier et al. [1]), combined with $\lambda = 1/n$ (Eq. 14), converges to the expected loss, but it provides good guarantees only for large training sample ($n \approx 105$). Note that the latter bound is not directly minimized by our "optimal posterior", as opposed to the one with $\lambda = 1/n$ (Eq. 13), for which we observe values between 5.8 (for $n=106$) and 6.4 (for $n=10$) not displayed on Figure 1c.

7

Conclusion

The first contribution of this paper is to bridge the concepts underlying the Bayesian and the PACBayesian approaches; under proper parameterization, the minimization of the PAC-Bayesian bound maximizes the marginal likelihood. This study motivates the second contribution of this paper, which is to prove PAC-Bayesian generalization bounds for regression with unbounded sub-gamma

loss functions, including the squared loss used in regression tasks. In this work, we studied model selection techniques. On a broader perspective, we would like to suggest that both Bayesian and PAC-Bayesian frameworks may have more to learn from each other than what has been done lately (even if other works paved the way [e.g., 6, 14, 30]). Predictors learned from the Bayes rule can benefit from strong PAC-Bayesian frequentist guarantees (under the i.i.d. assumption). Also, the rich Bayesian toolbox may be incorporated in PAC-Bayesian driven algorithms and risk bounding techniques. Acknowledgments We thank Gabriel Dub? and Maxime Tremblay for having proofread the paper and supplemental.

8

2 References

- [1] Pierre Alquier, James Ridgway, and Nicolas Chopin. On the properties of variational approximations of Gibbs posteriors. *JMLR*, 17(239):1?41, 2016.
- [2] Amiran Ambroladze, Emilio Parrado-Hern?ndez, and John Shawe-Taylor. Tighter PAC-Bayes bounds. In *NIPS*, 2006.
- [3] Arindam Banerjee. On Bayesian bounds. In *ICML*, pages 81?88, 2006.
- [4] Luc B?gin, Pascal Germain, Fran?ois Laviolette, and Jean-Fran?ois Roy. PAC-Bayesian theory for transductive learning. In *AISTATS*, pages 105?113, 2014.
- [5] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.
- [6] P. G. Bissiri, C. C. Holmes, and S. G. Walker. A general framework for updating belief distributions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2016.
- [7] St?phane Boucheron, G?bor Lugosi, and Pascal Massart. *Concentration inequalities : a nonasymptotic theory of independence*. Oxford university press, 2013. ISBN 978-0-19-953525-5.
- [8] Olivier Catoni. *PAC-Bayesian supervised classification: the thermodynamics of statistical learning*, volume 56. Inst. of Mathematical Statistic, 2007.
- [9] Arnak S. Dalalyan and Alexandre B. Tsybakov. Aggregation by exponential weighting, sharp PAC-Bayesian bounds and sparsity. *Machine Learning*, 72(1-2):39?61, 2008.
- [10] Pascal Germain, Alexandre Lacasse, Fran?ois Laviolette, and Mario Marchand. PAC-Bayesian learning of linear classifiers. In *ICML*, pages 353?360, 2009.
- [11] Pascal Germain, Alexandre Lacasse, Fran?ois Laviolette, Mario Marchand, and Jean-Fran?ois Roy. Risk bounds for the majority vote: From a PAC-Bayesian analysis to a learning algorithm. *JMLR*, 16, 2015.
- [12] Pascal Germain, Amaury Habrard, Fran?ois Laviolette, and Emilie Morvant. A new PAC-Bayesian perspective on domain adaptation. In *ICML*, pages 859?868, 2016.
- [13] Zoubin Ghahramani. Probabilistic machine learning and artificial intelligence. *Nature*, 521:452?459, 2015.
- [14] Peter Gr?nwald. The safe Bayesian - learning the learning rate via the mixability gap. In *ALT*, 2012.
- [15] Peter D. Gr?nwald and Nishant A. Mehta. Fast rates with unbounded losses. *CoRR*, abs/1605.00252, 2016.
- [16] Isabelle Guyon, Amir Saffari, Gideon Dror, and Gavin C. Cawley. Model selection: Beyond the Bayesian/frequentist divide. *JMLR*, 11:61?87, 2010.
- [17] Tamir Hazan, Subhransu Maji, Joseph Keshet, and

Tommi S. Jaakkola. Learning efficient random maximum a-posteriori predictors with non-decomposable loss functions. In NIPS, pages 1887–1895, 2013.

[18] William H. Jeffreys and James O. Berger. Ockham’s razor and Bayesian analysis. *American Scientist*, 1992.

[19] Alexandre Lacoste. Agnostic Bayes. PhD thesis, Université Laval, 2015.

[20] John Langford and John Shawe-Taylor. PAC-Bayes & margins. In NIPS, pages 423–430, 2002.

[21] Guy Lever, François Laviolette, and John Shawe-Taylor. Tighter PAC-Bayes bounds through distribution-independent priors. *Theor. Comput. Sci.*, 473:4–28, 2013.

[22] David J. C. MacKay. Bayesian interpolation. *Neural Computation*, 4(3):415–447, 1992.

[23] Andreas Maurer. A note on the PAC-Bayesian theorem. CoRR, cs.LG/0411099, 2004.

[24] David McAllester. Some PAC-Bayesian theorems. *Machine Learning*, 37(3):355–363, 1999.

[25] David McAllester. PAC-Bayesian stochastic model selection. *Machine Learning*, 51(1):5–21, 2003.

[26] David McAllester and Joseph Keshet. Generalization bounds and consistency for latent structural probit and ramp loss. In NIPS, pages 2205–2212, 2011.

[27] Asf Noy and Koby Crammer. Robust forward algorithms via PAC-Bayes and Laplace distributions. In AISTATS, 2014.

[28] Anastasia Pentina and Christoph H. Lampert. A PAC-Bayesian bound for lifelong learning. In ICML, 2014.

[29] Matthias Seeger. PAC-Bayesian generalization bounds for Gaussian processes. *JMLR*, 3:233–269, 2002.

[30] Matthias Seeger. Bayesian Gaussian Process Models: PAC-Bayesian Generalisation Error Bounds and Sparse Approximations. PhD thesis, University of Edinburgh, 2003.

[31] Yevgeny Seldin and Naftali Tishby. PAC-Bayesian analysis of co-clustering and beyond. *JMLR*, 11, 2010.

[32] Yevgeny Seldin, Peter Auer, François Laviolette, John Shawe-Taylor, and Ronald Ortner. PAC-Bayesian analysis of contextual bandits. In NIPS, pages 1683–1691, 2011.

[33] Yevgeny Seldin, François Laviolette, Nicolò Cesa-Bianchi, John Shawe-Taylor, and Peter Auer. PAC-Bayesian inequalities for martingales. In UAI, 2012.

[34] John Shawe-Taylor and Robert C. Williamson. A PAC analysis of a Bayesian estimator. In COLT, 1997.

[35] Ilya O. Tolstikhin and Yevgeny Seldin. PAC-Bayes-empirical-Bernstein inequality. In NIPS, 2013.

[36] Tong Zhang. Information-theoretic upper and lower bounds for statistical estimation. *IEEE Trans. Information Theory*, 52(4):1307–1321, 2006.