
PAC-Bayesian Theory Meets Bayesian Inference

Pascal Germain[†] Francis Bach[†] Alexandre Lacoste[‡] Simon Lacoste-Julien[†]

[†] INRIA Paris - École Normale Supérieure, `firstname.lastname@inria.fr`

[‡] Google, `allac@google.com`

Abstract

We exhibit a strong link between frequentist PAC-Bayesian risk bounds and the Bayesian marginal likelihood. That is, for the negative log-likelihood loss function, we show that the minimization of PAC-Bayesian generalization risk 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 Introduction

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 $1-\delta$, the generalization risk is at most ε 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 $(X, Y) = \{(x_i, y_i)\}_{i=1}^n \in (\mathcal{X} \times \mathcal{Y})^n$, that contains n input-output pairs. The main assumption of frequentist learning theories—including PAC-Bayes—is that (X, Y) is

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

randomly sampled from a data generating distribution that we denote \mathcal{D} . Thus, we denote $(X, Y) \sim \mathcal{D}^n$ the *i.i.d.* observation of n elements. From a frequentist perspective, we consider in this work loss functions $\ell : \mathcal{F} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, where \mathcal{F} is a (discrete or continuous) set of predictors $f : \mathcal{X} \rightarrow \mathcal{Y}$, and we write the empirical risk on the sample (X, Y) and the generalization error on distribution \mathcal{D} as

$$\widehat{\mathcal{L}}_{X,Y}^\ell(f) = \frac{1}{n} \sum_{i=1}^n \ell(f, x_i, y_i); \quad \mathcal{L}_\mathcal{D}^\ell(f) = \mathbf{E}_{(x,y) \sim \mathcal{D}} \ell(f, x, y).$$

The PAC-Bayesian theory [24, 25] studies an averaging of the above losses according to a *posterior* distribution $\hat{\rho}$ over \mathcal{F} . That is, it provides *probably approximately correct* generalization bounds on the (unknown) quantity $\mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_\mathcal{D}^\ell(f) = \mathbf{E}_{f \sim \hat{\rho}} \mathbf{E}_{(x,y) \sim \mathcal{D}} \ell(f, x, y)$, given the empirical estimate $\mathbf{E}_{f \sim \hat{\rho}} \widehat{\mathcal{L}}_{X,Y}^\ell(f)$ and some other parameters. Among these, most PAC-Bayesian theorems rely on the *Kullback-Leibler* divergence $\text{KL}(\hat{\rho} \parallel \pi) = \mathbf{E}_{f \sim \hat{\rho}} \ln[\hat{\rho}(f)/\pi(f)]$ between a *prior* distribution π over \mathcal{F} —specified before seeing the learning sample X, Y —and the posterior $\hat{\rho}$ —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 $\hat{\rho}$ over \mathcal{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 \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, a hypothesis set \mathcal{F} , a loss function $\ell' : \mathcal{F} \times \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$, a prior distribution π over \mathcal{F} , a real number $\delta \in (0, 1]$, and a real number $\beta > 0$, with probability at least $1 - \delta$ over the choice of $(X, Y) \sim \mathcal{D}^n$, we have*

$$\forall \hat{\rho} \text{ on } \mathcal{F} : \quad \mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_\mathcal{D}^{\ell'}(f) \leq \frac{1}{1 - e^{-\beta}} \left[1 - e^{-\beta \mathbf{E}_{f \sim \hat{\rho}} \widehat{\mathcal{L}}_{X,Y}^{\ell'}(f) - \frac{1}{n} (\text{KL}(\hat{\rho} \parallel \pi) + \ln \frac{1}{\delta})} \right]. \quad (1)$$

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 $\beta := b - a$ and with the *rescaled* loss function $\ell'(f, x, y) := (\ell(f, x, y) - a)/(b - a) \in [0, 1]$. After few arithmetic manipulations, we can rewrite Equation (1) as

$$\forall \hat{\rho} \text{ on } \mathcal{F} : \quad \mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_\mathcal{D}^\ell(f) \leq a + \frac{b-a}{1 - e^{-\beta}} \left[1 - \exp \left(-\mathbf{E}_{f \sim \hat{\rho}} \widehat{\mathcal{L}}_{X,Y}^\ell(f) + a - \frac{1}{n} (\text{KL}(\hat{\rho} \parallel \pi) + \ln \frac{1}{\delta}) \right) \right]. \quad (2)$$

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 π , X , Y , n , and δ , minimizing Equation (2) is equivalent to find the distribution $\hat{\rho}$ that minimizes

$$n \mathbf{E}_{f \sim \hat{\rho}} \widehat{\mathcal{L}}_{X,Y}^\ell(f) + \text{KL}(\hat{\rho} \parallel \pi). \quad (3)$$

It is well known [1, 8, 10, 21] that the *optimal Gibbs posterior* $\hat{\rho}^*$ is given by

$$\hat{\rho}^*(f) = \frac{1}{Z_{X,Y}} \pi(f) e^{-n \widehat{\mathcal{L}}_{X,Y}^\ell(f)}, \quad (4)$$

where $Z_{X,Y}$ is a normalization term. Notice that the constant β of Equation (1) is now absorbed in the loss function as the rescaling factor setting the trade-off between the expected empirical loss and $\text{KL}(\hat{\rho} \parallel \pi)$.

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.*, $p(y|x, \theta) \equiv f_\theta(y|x)$.² Using Bayes' rule, we obtain the posterior $p(\theta|X, Y)$:

$$p(\theta|X, Y) = \frac{p(\theta) p(Y|X, \theta)}{p(Y|X)} \propto p(\theta) p(Y|X, \theta), \quad (5)$$

where $p(Y|X, \theta) = \prod_{i=1}^n p(y_i|x_i, \theta)$ and $p(Y|X) = \int_\Theta p(\theta) p(Y|X, \theta) d\theta$.

²To 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.

To bridge the Bayesian approach with the PAC-Bayesian framework, we consider the *negative log-likelihood* loss function [3], denoted ℓ_{nl} and defined by

$$\ell_{\text{nl}}(f_\theta, x, y) \equiv -\ln p(y|x, \theta). \quad (6)$$

Then, we can relate the *empirical loss* $\widehat{\mathcal{L}}_{X,Y}^\ell$ of a predictor to its likelihood:

$$\widehat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{nl}}(\theta, x_i, y_i) = -\frac{1}{n} \sum_{i=1}^n \ln p(y_i|x_i, \theta) = -\frac{1}{n} \ln p(Y|X, \theta),$$

or, the other way around,

$$p(Y|X, \theta) = e^{-n \widehat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta)}. \quad (7)$$

Unfortunately, existing PAC-Bayesian theorems work with bounded loss functions or in very specific contexts [e.g., 9, 36], and ℓ_{nl} 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 $\ln \frac{1}{p(y|x, \theta)} \notin [a, b]$.

Now, using Equation (7) in the optimal posterior (Equation 4) simplifies to

$$\hat{\rho}^*(\theta) = \frac{\pi(\theta) e^{-n \widehat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta)}}{Z_{X,Y}} = \frac{p(\theta) p(Y|X, \theta)}{p(Y|X)} = p(\theta|X, Y), \quad (8)$$

where the normalization constant $Z_{X,Y}$ corresponds to the Bayesian *marginal likelihood*:

$$Z_{X,Y} \equiv p(Y|X) = \int_{\Theta} \pi(\theta) e^{-n \widehat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta)} d\theta. \quad (9)$$

This shows that the optimal PAC-Bayes posterior given by the generalization bound of Theorem 1 coincides with the Bayesian posterior, when one chooses ℓ_{nl} as loss function and $\beta := b-a$ (as in Equation 2). Moreover, using the posterior of Equation (8) inside Equation (3), we obtain

$$\begin{aligned} n \mathbf{E}_{\theta \sim \hat{\rho}^*} \widehat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta) + \text{KL}(\hat{\rho}^* \parallel \pi) & \quad (10) \\ &= n \int_{\Theta} \frac{\pi(\theta) e^{-n \widehat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta)}}{Z_{X,Y}} \widehat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta) d\theta + \int_{\Theta} \frac{\pi(\theta) e^{-n \widehat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta)}}{Z_{X,Y}} \ln \left[\frac{\pi(\theta) e^{-n \widehat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta)}}{\pi(\theta) Z_{X,Y}} \right] d\theta \\ &= \int_{\Theta} \frac{\pi(\theta) e^{-n \widehat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta)}}{Z_{X,Y}} \left[\ln \frac{1}{Z_{X,Y}} \right] d\theta = \frac{Z_{X,Y}}{Z_{X,Y}} \ln \frac{1}{Z_{X,Y}} = -\ln Z_{X,Y}. \end{aligned}$$

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 Grünwald [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 \mathcal{D} , a parameter set Θ , a prior distribution π over Θ , a $\delta \in (0, 1]$, if ℓ_{nl} lies in $[a, b]$, we have, with probability at least $1 - \delta$ over the choice of $(X, Y) \sim \mathcal{D}^n$,*

$$\mathbf{E}_{\theta \sim \hat{\rho}^*} \mathcal{L}_{\mathcal{D}}^{\ell_{\text{nl}}}(\theta) \leq a + \frac{b-a}{1-e^{-a}} \left[1 - e^a \sqrt[n]{Z_{X,Y} \delta} \right],$$

where $\hat{\rho}^*$ is the Gibbs optimal posterior (Eq. 8) and $Z_{X,Y}$ 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_{nl} , 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 \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, a hypothesis set \mathcal{F} , a loss function $\ell : \mathcal{F} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, a prior distribution π over \mathcal{F} , a $\delta \in (0, 1]$, and a real number $\lambda > 0$, with probability at least $1 - \delta$ over the choice of $(X, Y) \sim \mathcal{D}^n$, we have*

$$\forall \hat{\rho} \text{ on } \mathcal{F}: \quad \mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_{\mathcal{D}}^{\ell}(f) \leq \mathbf{E}_{f \sim \hat{\rho}} \hat{\mathcal{L}}_{X,Y}^{\ell}(f) + \frac{1}{\lambda} \left[\text{KL}(\hat{\rho} \parallel \pi) + \ln \frac{1}{\delta} + \Psi_{\ell, \pi, \mathcal{D}}(\lambda, n) \right], \quad (11)$$

$$\text{where } \Psi_{\ell, \pi, \mathcal{D}}(\lambda, n) = \ln \mathbf{E}_{f \sim \pi} \mathbf{E}_{X', Y' \sim \mathcal{D}^n} \exp \left[\lambda \left(\mathcal{L}_{\mathcal{D}}^{\ell}(f) - \hat{\mathcal{L}}_{X', Y'}^{\ell}(f) \right) \right]. \quad (12)$$

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

$$\forall \hat{\rho} \text{ on } \mathcal{F}: \quad \mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_{\mathcal{D}}^{\ell}(f) \leq \mathbf{E}_{f \sim \hat{\rho}} \hat{\mathcal{L}}_{X,Y}^{\ell}(f) + \frac{1}{n} [\text{KL}(\hat{\rho} \parallel \pi) + \ln \frac{1}{\delta}] + \frac{1}{2}(b-a)^2. \quad (13)$$

Note that the latter bound leads to the same trade-off as Theorem 1 (expressed by Equation 3). However, the choice $\lambda := n$ has the inconvenience that the bound value is at least $\frac{1}{2}(b-a)^2$, even at the limit $n \rightarrow \infty$. With $\lambda := \sqrt{n}$ the bound converges (a result similar to Equation (14) is also formulated by Pentina and Lampert [28]):

$$\forall \hat{\rho} \text{ on } \mathcal{F}: \quad \mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_{\mathcal{D}}^{\ell}(f) \leq \mathbf{E}_{f \sim \hat{\rho}} \hat{\mathcal{L}}_{X,Y}^{\ell}(f) + \frac{1}{\sqrt{n}} [\text{KL}(\hat{\rho} \parallel \pi) + \ln \frac{1}{\delta} + \frac{1}{2}(b-a)^2]. \quad (14)$$

Sub-Gaussian losses. In a regression context, it may be restrictive to consider strictly bounded loss functions. Therefore, we extend Theorem 3 to *sub-Gaussian* losses. We say that a loss function ℓ is sub-Gaussian with variance factor s^2 under a prior π and a data-distribution \mathcal{D} if it can be described by a sub-Gaussian random variable $V = \mathcal{L}_{\mathcal{D}}^{\ell}(f) - \ell(f, x, y)$, i.e., its moment generating function is upper bounded by the one of a normal distribution of variance s^2 (see Boucheron et al. [7, Section 2.3]):

$$\psi_V(\lambda) = \ln \mathbf{E} e^{\lambda V} = \ln \mathbf{E}_{f \sim \pi} \mathbf{E}_{(x,y) \sim \mathcal{D}} \exp [\lambda (\mathcal{L}_{\mathcal{D}}^{\ell}(f) - \ell(f, x, y))] \leq \frac{\lambda^2 s^2}{2}, \quad \forall \lambda \in \mathbb{R}. \quad (15)$$

The above sub-Gaussian assumption corresponds to the *Hoeffding assumption* of Alquier et al. [1], and allows to obtain the following result.

Corollary 4. *Given \mathcal{D} , \mathcal{F} , ℓ , π and δ defined in the statement of Theorem 3, if the loss is sub-Gaussian with variance factor s^2 , we have, with probability at least $1 - \delta$ over the choice of $(X, Y) \sim \mathcal{D}^n$,*

$$\forall \hat{\rho} \text{ on } \mathcal{F}: \quad \mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_{\mathcal{D}}^{\ell}(f) \leq \mathbf{E}_{f \sim \hat{\rho}} \hat{\mathcal{L}}_{X,Y}^{\ell}(f) + \frac{1}{n} [\text{KL}(\hat{\rho} \parallel \pi) + \ln \frac{1}{\delta}] + \frac{1}{2} s^2.$$

Proof. For $i = 1 \dots n$, we denote ℓ_i a i.i.d. realization of the random variable $\mathcal{L}_{\mathcal{D}}^{\ell}(f) - \ell(f, x, y)$.

$\Psi_{\ell, \pi, \mathcal{D}}(\lambda, n) = \ln \mathbf{E} \exp \left[\frac{\lambda}{n} \sum_{i=1}^n \ell_i \right] = \ln \prod_{i=1}^n \mathbf{E} \exp \left[\frac{\lambda}{n} \ell_i \right] = \sum_{i=1}^n \psi_{\ell_i} \left(\frac{\lambda}{n} \right) \leq n \frac{\lambda^2 s^2}{2n^2} = \frac{\lambda^2 s^2}{2n}$, where the inequality comes from the sub-Gaussian loss assumption (Equation 15). The result is then obtained from Theorem 3, with $\lambda := n$. \square

Sub-gamma losses. We say that an unbounded loss function ℓ is sub-gamma with a variance factor s^2 and scale parameter c , under a prior π and a data-distribution \mathcal{D} , if it can be described by a sub-gamma random variable V (see Boucheron et al. [7, Section 2.4]), that is

$$\psi_V(\lambda) \leq \frac{s^2}{c^2} (-\ln(1 - \lambda c) - \lambda c) \leq \frac{\lambda^2 s^2}{2(1 - c\lambda)}, \quad \forall \lambda \in (0, \frac{1}{c}). \quad (16)$$

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

Corollary 5. Given \mathcal{D} , \mathcal{F} , ℓ , π and δ defined in the statement of Theorem 3, if the loss is sub-gamma with variance factor s^2 and scale $c < 1$, we have, with probability at least $1 - \delta$ over $(X, Y) \sim \mathcal{D}^n$,

$$\forall \hat{\rho} \text{ on } \mathcal{F}: \quad \mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_{\mathcal{D}}^{\ell}(f) \leq \mathbf{E}_{f \sim \hat{\rho}} \hat{\mathcal{L}}_{X,Y}^{\ell}(f) + \frac{1}{n} [\text{KL}(\hat{\rho} \parallel \pi) + \ln \frac{1}{\delta}] + \frac{1}{2(1-c)} s^2. \quad (17)$$

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

$$\mathbf{E}_{\theta \sim \hat{\rho}^*} \mathcal{L}_{\mathcal{D}}^{\ell_{\text{nl}}}(\theta) \leq \frac{s^2}{2(1-c)} - \frac{1}{n} \ln(Z_{X,Y} \delta). \quad (18)$$

Proof. Following the same path as in the proof of Corollary 4 (with $\lambda := n$), we have

$$\Psi_{\ell, \pi, \mathcal{D}}(n, n) = \ln \mathbf{E} \exp \left[\sum_{i=1}^n \ell_i \right] = \ln \prod_{i=1}^n \mathbf{E} \exp [\ell_i] = \sum_{i=1}^n \psi_{\ell_i}(1) \leq n \frac{s^2}{2(1-c)} = \frac{n s^2}{2(1-c)},$$

where the inequality comes from the sub-gamma loss assumption, with $1 \in (0, \frac{1}{c})$. \square

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 $\mathcal{X} \times \mathcal{Y} \subset \mathbb{R}^d \times \mathbb{R}$, a family of linear predictors $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$, with $\mathbf{w} \in \mathbb{R}^d$, and a Gaussian prior $\mathcal{N}(\mathbf{0}, \sigma_{\pi}^2 \mathbf{I})$. Let us assume that the input examples are generated by $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \sigma_{\mathbf{x}}^2 \mathbf{I})$ with label $y = \mathbf{w}^* \cdot \mathbf{x} + \epsilon$, where $\mathbf{w}^* \in \mathbb{R}^d$ and $\epsilon \sim \mathcal{N}(0, \sigma_{\epsilon}^2)$ is a Gaussian noise. Under the squared loss function

$$\ell_{\text{sqf}}(\mathbf{w}, \mathbf{x}, y) = (\mathbf{w} \cdot \mathbf{x} - y)^2, \quad (19)$$

we show in Appendix A.4 that Corollary 5 is valid with $s^2 \geq 2 [\sigma_{\mathbf{x}}^2(\sigma_{\pi}^2 d + \|\mathbf{w}^*\|^2) + \sigma_{\epsilon}^2(1 - c)]$ and $c \geq 2\sigma_{\mathbf{x}}^2\sigma_{\pi}^2$. As expected, the bound degrades when the noise increases

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* $G_{\hat{\rho}}$. In order to predict the label of an example $x \in \mathcal{X}$, the Gibbs predictor first draws a hypothesis $h \in \mathcal{F}$ according to $\hat{\rho}$, 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 $y \in \{-1, 1\}$ and $h(x) \in [-1, 1]$, a common choice is to use the linear loss function $\ell'_{01}(h, x, y) = \frac{1}{2} - \frac{1}{2} y h(x)$. The Gibbs generalization loss is then given by $R_{\mathcal{D}}(G_{\hat{\rho}}) = \mathbf{E}_{(x,y) \sim \mathcal{D}} \mathbf{E}_{h \sim \hat{\rho}} \ell'_{01}(h, x, y)$. Many PAC-Bayesian works use $R_{\mathcal{D}}(G_{\hat{\rho}})$ as a surrogate loss to study the zero-one classification loss of the majority vote classifier $R_{\mathcal{D}}(B_{\hat{\rho}})$:

$$R_{\mathcal{D}}(B_{\hat{\rho}}) = \Pr_{(x,y) \sim \mathcal{D}} \left(y \mathbf{E}_{h \sim \hat{\rho}} h(x) < 0 \right) = \mathbf{E}_{(x,y) \sim \mathcal{D}} I \left[y \mathbf{E}_{h \sim \hat{\rho}} h(x) < 0 \right], \quad (20)$$

where $I[\cdot]$ being the indicator function. Given a distribution $\hat{\rho}$, an upper bound on the Gibbs risk is converted to an upper bound on the majority vote risk by $R_{\mathcal{D}}(B_{\hat{\rho}}) \leq 2R_{\mathcal{D}}(G_{\hat{\rho}})$ [20]. In some situations, this *factor of two* may be reached, i.e., $R_{\mathcal{D}}(B_{\hat{\rho}}) \simeq 2R_{\mathcal{D}}(G_{\hat{\rho}})$. In other situations, we may have $R_{\mathcal{D}}(B_{\hat{\rho}}) = 0$ even if $R_{\mathcal{D}}(G_{\hat{\rho}}) = \frac{1}{2} - \epsilon$ (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 $\mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_{\mathcal{D}}^{\ell}(f)$ for any (bounded, sub-Gaussian, or sub-gamma) losses. Of course, one may want to bound the regression loss of the averaged regressor $F_{\hat{\rho}}(x) = \mathbf{E}_{f \sim \hat{\rho}} f(x)$. In this case, if the loss function ℓ is convex (as the squared loss), Jensen's inequality gives $\mathcal{L}_{\mathcal{D}}^{\ell}(F_{\hat{\rho}}) \leq \mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_{\mathcal{D}}^{\ell}(f)$. 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.

³It is noteworthy that the best PAC-Bayesian empirical bound values are so far obtained by considering a majority vote of linear classifiers, where the prior and posterior are Gaussian [2, 10, 20], similarly to the Bayesian linear regression analyzed in Section 6.

5 Analysis of Model Selection

We consider L distinct models $\{\mathcal{M}_i\}_{i=1}^L$, each one defined by a set of parameters Θ_i . The PAC-Bayesian theorems naturally suggest selecting the model that is best adapted for the given task by evaluating the bound for each model $\{\mathcal{M}_i\}_{i=1}^L$ 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 posteriors—one for each model—given by Equation (8),

$$p(\theta|X, Y, \mathcal{M}_i) \equiv \hat{\rho}_i^*(\theta) = \frac{1}{Z_{X,Y,i}} \pi_i(\theta) e^{-n \hat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta)}, \text{ for } \theta \in \Theta_i, \quad (21)$$

the Bayesian Occam’s razor criteria [18, 22] chooses the one with the higher *model evidence*

$$p(Y|X, \mathcal{M}_i) \equiv Z_{X,Y,i} = \int_{\Theta_i} \pi_i(\theta) e^{-n \hat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta)} d\theta. \quad (22)$$

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 ℓ_{nl} and Equation (10). From the union bound (*a.k.a.* Bonferroni inequality), it is mandatory to compute each bound with a confidence parameter of δ/L , to ensure that the final conclusion is valid with probability at least $1-\delta$.

Corollary 6. *Given a data distribution \mathcal{D} , a family of model parameters $\{\Theta_i\}_{i=1}^L$ and associated priors $\{\pi_i\}_{i=1}^L$ —where π_i is defined over Θ_i —, a $\delta \in (0, 1]$, if the loss is sub-gamma with parameters s^2 and $c < 1$, then, with probability at least $1 - \delta$ over $(X, Y) \sim \mathcal{D}^n$,*

$$\forall i \in \{1, \dots, L\} : \quad \mathbf{E}_{\theta \sim \hat{\rho}_i^*} \mathcal{L}_{\mathcal{D}}^{\ell_{\text{nl}}}(\theta) \leq \frac{1}{2(1-c)} s^2 - \frac{1}{n} \ln \left(Z_{X,Y,i} \frac{\delta}{L} \right).$$

where $\hat{\rho}_i^*$ 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* $p(\eta)$ over the set of hyperparameters \mathcal{H} . Then, the prior $p(\theta|\eta)$ can be conditioned on a choice of hyperparameter η . The Bayes rule of Equation (5) becomes $p(\theta, \eta|X, Y) = \frac{p(\eta) p(\theta|\eta) p(Y|X, \theta)}{p(Y|X)}$.

Under the negative log-likelihood loss function, we can rewrite the results of Corollary 5 as a generalization bound on $\mathbf{E}_{\eta \sim \hat{\rho}_0} \mathbf{E}_{\theta \sim \hat{\rho}_\eta^*} \mathcal{L}_{\mathcal{D}}^{\ell_{\text{nl}}}(\theta)$, where $\hat{\rho}_0(\eta) \propto \pi_0(\eta) Z_{X,Y,\eta}$ is the hyperposterior on \mathcal{H} and π_0 the hyperprior. Indeed, Equation (18) becomes

$$\mathbf{E}_{\theta \sim \hat{\rho}^*} \mathcal{L}_{\mathcal{D}}^{\ell_{\text{nl}}}(\theta) = \mathbf{E}_{\eta \sim \hat{\rho}_0} \mathbf{E}_{\theta \sim \hat{\rho}_\eta^*} \mathcal{L}_{\mathcal{D}}^{\ell_{\text{nl}}}(\theta) \leq \frac{1}{2(1-c)} s^2 - \frac{1}{n} \ln \left(\mathbf{E}_{\eta \sim \pi_0} Z_{X,Y,\eta} \frac{\delta}{L} \right). \quad (23)$$

To relate to the bound obtained in Corollary 6, we consider the case of a discrete hyperparameter set $\mathcal{H} = \{\eta_i\}_{i=1}^L$, with a uniform prior $\pi_0(\eta_i) = \frac{1}{L}$ (from now on, we regard each hyperparameter η_i as the specification of a model Θ_i). Then, Equation (23) becomes

$$\mathbf{E}_{\theta \sim \hat{\rho}^*} \mathcal{L}_{\mathcal{D}}^{\ell_{\text{nl}}}(\theta) = \mathbf{E}_{\eta \sim \hat{\rho}_0} \mathbf{E}_{\theta \sim \hat{\rho}_\eta^*} \mathcal{L}_{\mathcal{D}}^{\ell_{\text{nl}}}(\theta) \leq \frac{1}{2(1-c)} s^2 - \frac{1}{n} \ln \left(\sum_{i=1}^L Z_{X,Y,\eta_i} \frac{\delta}{L} \right).$$

This bound is now a function of $\sum_{i=1}^L Z_{X,Y,\eta_i}$ instead of $\max_i Z_{X,Y,\eta_i}$ 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 $\eta^* \in \mathcal{H}$, the hierarchical representation is equivalent to choosing a deterministic hyperposterior, satisfying $\hat{\rho}_0(\eta^*) = 1$ and 0 for every other values. We then have

$$\text{KL}(\hat{\rho}||\pi) = \text{KL}(\hat{\rho}_0||\pi_0) + \mathbf{E}_{\eta \sim \hat{\rho}_0} \text{KL}(\hat{\rho}_\eta||\pi_\eta) = \ln(L) + \text{KL}(\hat{\rho}_{\eta^*}||\pi_{\eta^*}).$$

With the optimal posterior for the selected η^* , we have

$$\begin{aligned} n \mathbf{E}_{\theta \sim \hat{\rho}} \hat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta) + \text{KL}(\hat{\rho}||\pi) &= n \mathbf{E}_{\theta \sim \hat{\rho}_{\eta^*}^*} \hat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\theta) + \text{KL}(\hat{\rho}_{\eta^*}^*||\pi_{\eta^*}) + \ln(L) \\ &= -\ln(Z_{X,Y,\eta^*}) + \ln(L) = -\ln \left(\frac{Z_{X,Y,\eta^*}}{L} \right). \end{aligned}$$

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.

6 Linear Regression

In this section, we perform *Bayesian linear regression* using the parameterization of Bishop [5]. The output space is $\mathcal{Y} := \mathbb{R}$ and, for an arbitrary input space \mathcal{X} , we use a mapping function $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$.

The model. Given $(x, y) \in \mathcal{X} \times \mathcal{Y}$ and model parameters $\theta := \langle \mathbf{w}, \sigma \rangle \in \mathbb{R}^d \times \mathbb{R}^+$, we consider the likelihood $p(y|x, \langle \mathbf{w}, \sigma \rangle) = \mathcal{N}(y|\mathbf{w} \cdot \phi(\mathbf{x}), \sigma^2)$. Thus, the negative log-likelihood loss is

$$\ell_{\text{nl}}(\langle \mathbf{w}, \sigma \rangle, x, y) = -\ln p(y|x, \langle \mathbf{w}, \sigma \rangle) = \frac{1}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} (y - \mathbf{w} \cdot \phi(x))^2. \quad (24)$$

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 $\mathbf{0}$ and variance σ_π^2 : $p(\mathbf{w}|\sigma_\pi) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \sigma_\pi^2 \mathbf{I})$. For the sake of simplicity, we consider fixed parameters σ^2 and σ_π^2 . The Gibbs optimal posterior (see Equation 8) is then given by

$$\hat{\rho}^*(\mathbf{w}) \equiv p(\mathbf{w}|X, Y, \sigma, \sigma_\pi) = \frac{p(\mathbf{w}|\sigma_\pi)p(Y|X, \mathbf{w}, \sigma)}{p(Y|X, \sigma, \sigma_\pi)} = \mathcal{N}(\mathbf{w}|\hat{\mathbf{w}}, A^{-1}), \quad (25)$$

where $A := \frac{1}{\sigma^2} \Phi^T \Phi + \frac{1}{\sigma_\pi^2} \mathbf{I}$; $\hat{\mathbf{w}} := \frac{1}{\sigma^2} A^{-1} \Phi^T \mathbf{y}$; Φ is a $n \times d$ matrix such that the i^{th} line is $\phi(x_i)$; $\mathbf{y} := [y_1, \dots, y_n]$ is the labels-vector; and the negative log marginal likelihood is

$$\begin{aligned} -\ln p(Y|X, \sigma, \sigma_\pi) &= \frac{1}{2\sigma^2} \|\mathbf{y} - \Phi \hat{\mathbf{w}}\|^2 + \frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma_\pi^2} \|\hat{\mathbf{w}}\|^2 + \frac{1}{2} \log |A| + d \ln \sigma_\pi \\ &= \underbrace{n \hat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\hat{\mathbf{w}}) + \frac{1}{2\sigma^2} \text{tr}(\Phi^T \Phi A^{-1})}_{n \mathbf{E}_{\mathbf{w} \sim \hat{\rho}^*} \hat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\mathbf{w})} + \underbrace{\frac{1}{2\sigma_\pi^2} \text{tr}(A^{-1}) - \frac{d}{2} + \frac{1}{2\sigma_\pi^2} \|\hat{\mathbf{w}}\|^2 + \frac{1}{2} \log |A| + d \ln \sigma_\pi}_{\text{KL}(\mathcal{N}(\hat{\mathbf{w}}, A^{-1}) \parallel \mathcal{N}(\mathbf{0}, \sigma_\pi^2 \mathbf{I}))}. \end{aligned}$$

To obtain the second equality, we substitute $\frac{1}{2\sigma^2} \|\mathbf{y} - \Phi \hat{\mathbf{w}}\|^2 + \frac{n}{2} \ln(2\pi\sigma^2) = n \hat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\hat{\mathbf{w}})$ and insert $\frac{1}{2\sigma^2} \text{tr}(\Phi^T \Phi A^{-1}) + \frac{1}{2\sigma_\pi^2} \text{tr}(A^{-1}) = \frac{1}{2} \text{tr}(\frac{1}{\sigma^2} \Phi^T \Phi A^{-1} + \frac{1}{\sigma_\pi^2} A^{-1}) = \frac{1}{2} \text{tr}(A^{-1} A) = \frac{d}{2}$. This exhibits how the Bayesian regression optimization problem is related to the minimization of a PAC-Bayesian bound, expressed by a trade-off between $\mathbf{E}_{\mathbf{w} \sim \hat{\rho}^*} \hat{\mathcal{L}}_{X,Y}^{\ell_{\text{nl}}}(\mathbf{w})$ and $\text{KL}(\mathcal{N}(\hat{\mathbf{w}}, A^{-1}) \parallel \mathcal{N}(\mathbf{0}, \sigma_\pi^2 \mathbf{I}))$. See Appendix A.5 for detailed calculations.

Model selection experiment. To produce Figures 1a and 1b, we reimplemented the toy experiment of Bishop [5, Section 3.5.1]. That is, we generated a learning sample of 15 data points according to $y = \sin(x) + \epsilon$, where x is uniformly sampled in the interval $[0, 2\pi]$ and $\epsilon \sim \mathcal{N}(0, \frac{1}{4})$ is a Gaussian noise. We then learn seven different polynomial models applying Equation (25). More precisely, for a polynomial model of degree d , we map input $x \in \mathbb{R}$ to a vector $\phi(x) = [1, x^1, x^2, \dots, x^d] \in \mathbb{R}^{d+1}$, and we fix parameters $\sigma_\pi^2 = \frac{1}{0.005}$ and $\sigma^2 = \frac{1}{2}$. Figure 1a illustrates the seven learned models. Figure 1b shows the negative log marginal likelihood computed for each polynomial model, and 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 $\mathbf{x} \in \mathbb{R}^{20}$ is generated by a Gaussian $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The associated output $y \in \mathbb{R}$ is given by $y = \mathbf{w}^* \cdot \mathbf{x} + \epsilon$, with $\|\mathbf{w}^*\| = \frac{1}{2}$, $\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2)$, and $\sigma_\epsilon^2 = \frac{1}{9}$. We perform Bayesian linear regression in the input space, i.e., $\phi(\mathbf{x}) = \mathbf{x}$, fixing $\sigma_\pi^2 = \frac{1}{100}$ and $\sigma^2 = 2$. That is, we compute the posterior of Equation (25) for training samples of sizes from 10 to 10^6 . For each learned model, we compute the empirical negative log-likelihood loss of Equation (24), and the three PAC-Bayes bounds, with confidence parameter of $\delta = \frac{1}{20}$. Note that this loss function is an affine transformation of the squared loss studied in Section 4 (Equation 19), i.e., $\ell_{\text{nl}}(\langle \mathbf{w}, \sigma \rangle, \mathbf{x}, y) = \frac{1}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \ell_{\text{sqr}}(\mathbf{w}, \mathbf{x}, y)$. It turns out that ℓ_{nl} is sub-gamma with parameters $s^2 \geq \frac{1}{\sigma^2} [\sigma_\pi^2 d + \|\mathbf{w}^*\|^2] + \sigma_\epsilon^2 (1 - c)$ and $c \geq \frac{1}{\sigma^2} (\sigma_\pi^2 \sigma^2)$, as shown in Appendix A.6. The bounds of Corollary 5 are computed using the above mentioned values of $\|\mathbf{w}^*\|, d, \sigma, \sigma_\pi, \sigma_\epsilon$, leading

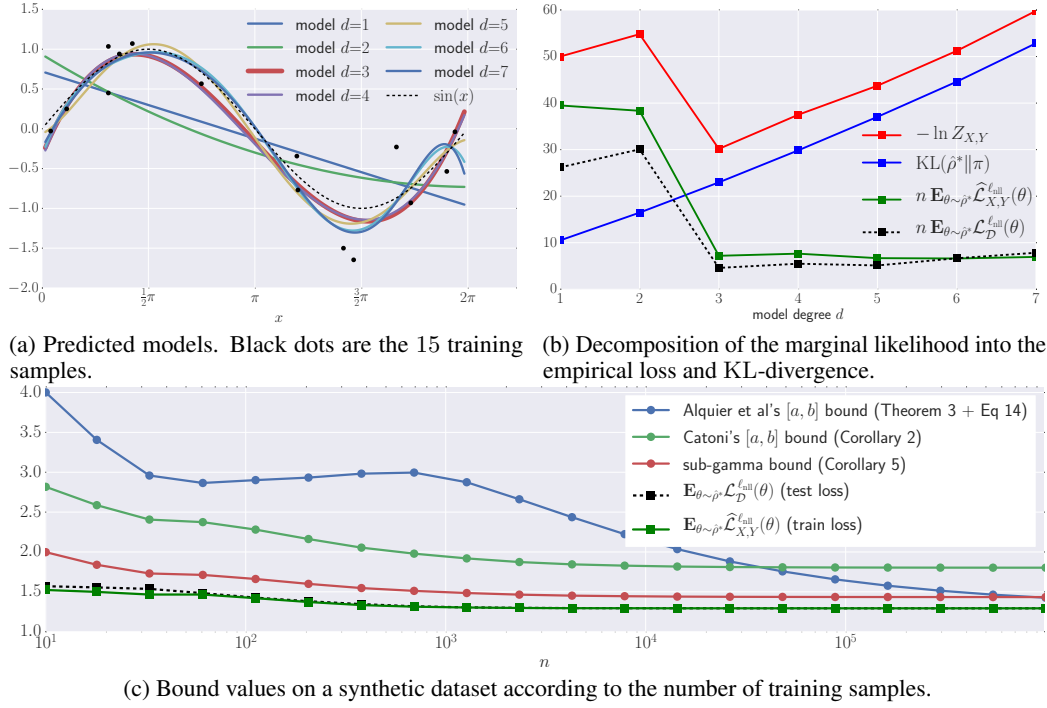


Figure 1: Model selection experiment (a-b); and comparison of bounds values (c).

to $s^2 \simeq 0.280$ and $c \simeq 0.005$. 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) sub-gamma losses—gives tighter guarantees than the two results for $[a, b]$ -bounded losses (up to $n=10^6$). However, our new bound always maintains a gap of $\frac{1}{2(1-c)}s^2$ 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 $\lambda = 1/\sqrt{n}$ (Eq. 14), converges to the expected loss, but it provides good guarantees only for large training sample ($n \gtrsim 10^5$). 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 $n=10^6$) 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 PAC-Bayesian 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 Dubé and Maxime Tremblay for having proofread the paper and supplemental.

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-Hernández, 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 Bégin, Pascal Germain, François Laviolette, and Jean-François 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] Stéphane Boucheron, Gábor 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, François Laviolette, and Mario Marchand. PAC-Bayesian learning of linear classifiers. In *ICML*, pages 353–360, 2009.
- [11] Pascal Germain, Alexandre Lacasse, François Laviolette, Mario Marchand, and Jean-François Roy. Risk bounds for the majority vote: From a PAC-Bayesian analysis to a learning algorithm. *JMLR*, 16, 2015.
- [12] Pascal Germain, Amaury Habrard, François 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 Grünwald. The safe Bayesian - learning the learning rate via the mixability gap. In *ALT*, 2012.
- [15] Peter D. Grünwald 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, Subhansu 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-dependent 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.

A Supplementary Material

A.1 Related Work

In this section, we discuss briefly other works containing (more or less indirect) links between Bayesian inference and PAC-Bayesian theory, and explain how they relate to the current paper.

Seeger (2002, 2003) [29, 30]. Soon after the initial work of McAllester [24, 25], Seeger shows how to apply the PAC-Bayesian theorems to bound the generalization error of Gaussian Processes in a classification context. By building upon the PAC-Bayesian theorem initially appearing in Langford and Seeger [44]—where the divergence between the training error and the generalization one is given by the Kullback-Leibler divergence between two Bernoulli distributions—it achieves very tight generalization bounds.⁴ Also, the thesis of Seeger [30, Section 3.2] foresees this by noticing that “the log marginal likelihood incorporates a *similar trade-off* as the PAC-Bayesian theorem”, but using another variant of the PAC-Bayes bound and in the context of classification.

Banerjee (2006) [3]. This paper shows similarities between the early PAC-Bayesian results (McAllester [25], Langford and Seeger [44]), and the *Bayesian log-loss bound* (Freund and Schapire [38], Kakade and Ng [42]). This is done by highlighting that the proof of all these results are strongly relying on the same *compression lemma* [3, Lemma 1], which is equivalent to our *change of measure* used in the proof of Theorem 3 (see forthcoming Equation 26). Note that the loss studied in the Bayesian part of Banerjee [3] is the negative log-likelihood of Equation (6). Also, as in Equation (10), the *Bayesian log-loss bound* contains the Kullback-Leibler divergence between the prior and the posterior. However, the latter result is not a generalization bound, but a bound on the training loss that is obtained by computing a surrogate training loss in the specific context of online learning. Moreover, the marginal likelihood and the model selection techniques are not addressed in Banerjee [3].

Zhang (2006) [36]. This paper presents a family of information theoretical bounds for *randomized estimators* that have a lot in common with PAC-Bayesian results (although the bounded quantity is not directly the generalization error). Minimizing these bounds leads to the same optimal Gibbs posterior of Equation (4). The author noted that using the negative log-likelihood (Equation 6) leads to the Bayesian posterior, but made no connection with the marginal likelihood.

Grünwald (2012) [14]. This paper proposes the *Safe Bayesian* algorithm, which selects a proper Bayesian *learning rate* — that is analogous to the parameter β of our Equation (1), and the parameter λ of our Equation (11) — in the context of *misspecified models*.⁵ The standard Bayesian inference method is obtained with a fixed learning rate, corresponding to the case $\lambda := n$ (that is the case we focus on the current paper, see Corollaries 4 and 5). The analysis of Grünwald [14] relies both on the Minimum Description Length principle [41] and PAC-Bayesian theory. Building upon the work of Zhang [36] discussed above, they formulate the result that we presented as Equation (10), linking the marginal likelihood to the inherent PAC-Bayesian trade-off. However, they do not compute explicit bounds on the generalization loss, which required us to take into account the complexity term of Equation (12).

Lacoste (2015) [19]. In a binary classification context, it is shown that the parameter β of Theorem 1 can be interpreted as a Bernoulli label noise model from a Bayesian likelihood standpoint. For more details, we refer the reader to Section 2.2 of this thesis.

Bissiri et al. (2016) [6]. This recent work studies Bayesian inference through the lens of loss functions. When the loss function is the negative log-likelihood (Equation 6), the approach of Bissiri et al. [6] coincides with the Bayesian update rule. As mentioned by the authors, there is some connection between their framework and the PAC-Bayesian one, but “the motivation and construction are very different.”

⁴The PAC-Bayesian results for Gaussian processes are summarized in Rasmussen and Williams [47, Section 7.4]

⁵The empirical model selection capabilities of the *Safe Bayesian* algorithm has been further studied in Grünwald and van Ommen [40].

Other references. See also Grünwald and Langford [39], Lacoste-Julien et al. [43], Meir and Zhang [45], Ng and Jordan [46], Rousseau [48] for other studies drawing links between frequentist statistics and Bayesian inference, but outside the PAC-Bayesian framework.

A.2 Proof of Theorem 3

Recall that Theorem 3 originally comes from Alquier et al. [1, Theorem 4.1]. We present below a different proof that follows the key steps of the very general PAC-Bayesian theorem presented in Bégin et al. [37, Theorem 4].

Proof of Theorem 3. The Donsker-Varadhan's change of measure states that, for any measurable function $\phi : \mathcal{F} \rightarrow \mathbb{R}$, we have

$$\mathbf{E}_{f \sim \hat{\rho}} \phi(f) \leq \text{KL}(\hat{\rho} \parallel \pi) + \ln \left(\mathbf{E}_{f \sim \pi} e^{\phi(f)} \right). \quad (26)$$

Thus, with $\phi(f) := \lambda(\mathcal{L}_{\mathcal{D}}^{\ell}(f) - \widehat{\mathcal{L}}_{X,Y}^{\ell}(f))$, we obtain $\forall \hat{\rho}$ on \mathcal{F} :

$$\begin{aligned} \lambda \left(\mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_{\mathcal{D}}^{\ell}(f) - \mathbf{E}_{f \sim \hat{\rho}} \widehat{\mathcal{L}}_{X,Y}^{\ell}(f) \right) &= \mathbf{E}_{f \sim \hat{\rho}} \lambda \left(\mathcal{L}_{\mathcal{D}}^{\ell}(f) - \widehat{\mathcal{L}}_{X,Y}^{\ell}(f) \right) \\ &\leq \text{KL}(\hat{\rho} \parallel \pi) + \ln \left(\mathbf{E}_{f \sim \pi} e^{\lambda(\mathcal{L}_{\mathcal{D}}^{\ell}(f) - \widehat{\mathcal{L}}_{X,Y}^{\ell}(f))} \right). \end{aligned}$$

Now, we apply Markov's inequality on the random variable $\zeta_{\pi}(X, Y) := \mathbf{E}_{f \sim \pi} e^{\lambda(\mathcal{L}_{\mathcal{D}}^{\ell}(f) - \widehat{\mathcal{L}}_{X,Y}^{\ell}(f))}$:

$$\Pr_{X,Y \sim \mathcal{D}^n} \left(\zeta_{\pi}(X, Y) \leq \frac{1}{\delta} \mathbf{E}_{X',Y' \sim \mathcal{D}^n} \zeta_{\pi}(X', Y') \right) \geq 1 - \delta.$$

This implies that with probability at least $1 - \delta$ over the choice of $X, Y \sim \mathcal{D}^n$, we have $\forall \hat{\rho}$ on \mathcal{F} :

$$\mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_{\mathcal{D}}^{\ell}(f) \leq \mathbf{E}_{f \sim \hat{\rho}} \widehat{\mathcal{L}}_{X,Y}^{\ell}(f) + \frac{1}{\lambda} \left[\text{KL}(\hat{\rho} \parallel \pi) + \ln \frac{\mathbf{E}_{X',Y' \sim \mathcal{D}^n} \zeta_{\pi}(X', Y')}{\delta} \right].$$

□

A.3 Proof of Equations (13) and (14)

Proof. Given a loss function $\ell : \mathcal{F} \times \mathcal{X} \times \mathcal{Y}$, and a fixed predictor $f \in \mathcal{F}$, we consider the random experiment of sampling $(x, y) \in \mathcal{D}$. We denote ℓ_i a realization of the random variable $\mathcal{L}_{\mathcal{D}}^{\ell}(f) - \ell(f, x, y)$, for $i = 1 \dots n$. Each ℓ_i is *i.i.d.*, zero mean, and bounded by $a - b$ and $b - a$, as $\ell(f, x, y) \in [a, b]$. Thus,

$$\begin{aligned} \mathbf{E}_{X',Y' \sim \mathcal{D}^n} \exp \left[\lambda \left(\mathcal{L}_{\mathcal{D}}^{\ell}(f) - \widehat{\mathcal{L}}_{X',Y'}^{\ell}(f) \right) \right] &= \mathbf{E} \exp \left[\frac{\lambda}{n} \sum_{i=1}^n \ell_i \right] \\ &= \prod_{i=1}^n \mathbf{E} \exp \left[\frac{\lambda}{n} \ell_i \right] \\ &\leq \prod_{i=1}^n \exp \left[\frac{\lambda^2 (a - b - (b - a))^2}{8n^2} \right] \\ &= \prod_{i=1}^n \exp \left[\frac{\lambda^2 (b - a)^2}{2n^2} \right] \\ &= \exp \left[\frac{\lambda^2 (b - a)^2}{2n} \right], \end{aligned}$$

where the inequality comes from Hoeffding's lemma.

With $\lambda := n$, Equation (11) becomes Equation (13) :

$$\begin{aligned}\mathbf{E}_{f \sim \hat{\rho}} \mathcal{L}_{\mathcal{D}}^{\ell}(f) &\leq \mathbf{E}_{f \sim \hat{\rho}} \hat{\mathcal{L}}_{X,Y}^{\ell}(f) + \frac{1}{n} \left[\text{KL}(\hat{\rho} \parallel \pi) + \ln \frac{1}{\delta} + \frac{n^2(b-a)^2}{2n} \right] \\ &= \mathbf{E}_{f \sim \hat{\rho}} \hat{\mathcal{L}}_{X,Y}^{\ell}(f) + \frac{1}{n} \left[\text{KL}(\hat{\rho} \parallel \pi) + \ln \frac{1}{\delta} \right] + \frac{1}{2}(b-a)^2.\end{aligned}$$

Similarly, with $\lambda := \sqrt{n}$, Equation (11) becomes Equation (14) . \square

A.4 Study of the Squared Loss

We consider a regression problem where $\mathcal{X} \times \mathcal{Y} \subset \mathbb{R}^d \times \mathbb{R}$, a family of linear predictors $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$, with $\mathbf{w} \in \mathbb{R}^d$, and a Gaussian prior $\mathcal{N}(\mathbf{0}, \sigma_{\pi}^2 \mathbf{I})$. Let us assume that the input examples are generated according to $\mathcal{N}(\mathbf{0}, \sigma_{\mathbf{x}}^2 \mathbf{I})$ and $\epsilon \sim \mathcal{N}(0, \sigma_{\epsilon}^2)$ is a Gaussian noise.

We study the squared loss $\ell_{\text{sqr}}(\mathbf{w}, \mathbf{x}, y) = (\mathbf{w} \cdot \mathbf{x} - y)^2$ such that:

- $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma_{\pi}^2 \mathbf{I})$ is given by the prior π ,
- $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \sigma_{\mathbf{x}}^2 \mathbf{I})$ (and $\mathbf{x} \in \mathbb{R}^d$),
- $y = \mathbf{w}^* \cdot \mathbf{x} + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma_{\epsilon}^2)$, corresponds to the labeling function.
Thus $y|\mathbf{x} \sim \mathcal{N}(\mathbf{x} \cdot \mathbf{w}^*, \sigma_{\epsilon}^2)$.

Let us consider the random variable $v = [\mathbf{E}_{\mathbf{x}} \mathbf{E}_{y|\mathbf{x}} \ell_{\text{sqr}}(\mathbf{w}, \mathbf{x}, y)] - \ell_{\text{sqr}}(\mathbf{w}, \mathbf{x}, y)$. To show that v is a sub-gamma random variable, we will find values of c and s such that the criterion of Equation (16) is fulfilled, *i.e.*,

$$\psi_v(\lambda) = \ln \mathbf{E} e^{\lambda v} \leq \frac{\lambda^2 s^2}{2(1-c\lambda)}, \quad \forall \lambda \in (0, \frac{1}{c}).$$

We have,

$$\begin{aligned}\psi_v(\lambda) &= \ln \mathbf{E}_{\mathbf{x}} \mathbf{E}_{y|\mathbf{x}} \mathbf{E}_{\mathbf{w}} \exp \left(\lambda [\mathbf{E}_{\mathbf{x}} \mathbf{E}_{y|\mathbf{x}} (y - \mathbf{w} \cdot \mathbf{x})^2] - \lambda (y - \mathbf{w} \cdot \mathbf{x})^2 \right) \\ &\leq \ln \mathbf{E}_{\mathbf{w}} \exp \left(\lambda \mathbf{E}_{\mathbf{x}} \mathbf{E}_{y|\mathbf{x}} (y - \mathbf{w} \cdot \mathbf{x})^2 \right) \\ &= \ln \mathbf{E}_{\mathbf{w}} \exp \left(\lambda \mathbf{E}_{\mathbf{x}} [\mathbf{x} \cdot (\mathbf{w}^* - \mathbf{w})]^2 + \lambda \sigma_{\epsilon}^2 \right) \\ &= \ln \mathbf{E}_{\mathbf{w}} \exp \left(\lambda \sigma_{\mathbf{x}}^2 \|\mathbf{w}^* - \mathbf{w}\|^2 + \lambda \sigma_{\epsilon}^2 \right) \\ &= \ln \frac{1}{(1 - 2\lambda \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2)^{\frac{d}{2}}} \exp \left(\frac{\lambda \sigma_{\mathbf{x}}^2 \|\mathbf{w}^*\|^2}{1 - 2\lambda \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2} + \lambda \sigma_{\epsilon}^2 \right) \\ &= -\frac{d}{2} \ln(1 - 2\lambda \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2) + \frac{\lambda \sigma_{\mathbf{x}}^2 \|\mathbf{w}^*\|^2}{1 - 2\lambda \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2} + \lambda \sigma_{\epsilon}^2 \\ &\leq \frac{\lambda \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2 d}{1 - 2\lambda \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2} + \frac{\lambda \sigma_{\mathbf{x}}^2 \|\mathbf{w}^*\|^2}{1 - 2\lambda \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2} + \lambda \sigma_{\epsilon}^2 \\ &= \frac{\lambda (\sigma_{\pi}^2 \sigma_{\mathbf{x}}^2 d + \sigma_{\mathbf{x}}^2 \|\mathbf{w}^*\|^2 + (1 - 2\lambda \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2) \sigma_{\epsilon}^2)}{1 - 2\lambda \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2} \\ &= \frac{\lambda^2 s^2}{2(1 - \lambda c)},\end{aligned}$$

with $s^2 = \frac{2}{\lambda} [\sigma_{\mathbf{x}}^2 (\sigma_{\pi}^2 d + \|\mathbf{w}^*\|^2) + \sigma_{\epsilon}^2 (1 - \lambda c)]$ and $c = 2\sigma_{\mathbf{x}}^2 \sigma_{\pi}^2$.

Recall that Corollary 5 is obtained with $\lambda := 1$.

A.5 Linear Regression : Detailed calculations

Recall that, from Equation (25), the Gibbs optimal posterior of the described model is given by

$$p(\mathbf{w} | X, Y, \sigma, \sigma_{\pi}) = \mathcal{N}(\mathbf{w} | \hat{\mathbf{w}}, A^{-1}),$$

with $A := \frac{1}{\sigma^2} \Phi^T \Phi + \frac{1}{\sigma_\pi^2} \mathbf{I}$; $\hat{\mathbf{w}} := \frac{1}{\sigma^2} A^{-1} \Phi^T \mathbf{y}$; Φ is a $n \times d$ matrix such that the i^{th} line is $\phi(x_i)$; $\mathbf{y} := [y_1, \dots, y_n]$ is the labels-vector. For the complete derivation leading to this posterior distribution, see Bishop [5, Section 3.3] or Rasmussen and Williams [47, Section 2.1.1].

Marginal likelihood. We decompose of the marginal likelihood into the PAC-Bayesian trade-off:

$$\begin{aligned} & -\ln p(Y|X, \sigma, \sigma_\pi) \\ &= \frac{1}{2\sigma^2} \|\mathbf{y} - \Phi \hat{\mathbf{w}}\|^2 + \frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma_\pi^2} \|\hat{\mathbf{w}}\|^2 + \frac{1}{2} \log |A| + d \ln \sigma_\pi \quad (\dagger) \\ &= \underbrace{n \hat{\mathcal{L}}_{X,Y}^{\text{nl}}(\hat{\mathbf{w}}) + \frac{1}{2\sigma^2} \text{tr}(\Phi^T \Phi A^{-1})}_{n \mathbf{E}_{\mathbf{w} \sim \hat{\rho}^*} \hat{\mathcal{L}}_{X,Y}^{\text{nl}}(\mathbf{w})} + \underbrace{\frac{1}{2\sigma_\pi^2} \text{tr}(A^{-1}) - \frac{d}{2} + \frac{1}{2\sigma_\pi^2} \|\hat{\mathbf{w}}\|^2 + \frac{1}{2} \log |A| + d \ln \sigma_\pi}_{\text{KL}(\mathcal{N}(\hat{\mathbf{w}}, A^{-1}) \parallel \mathcal{N}(\mathbf{0}, \sigma_\pi^2 \mathbf{I}))} \quad (\star) \end{aligned}$$

Line (\dagger) corresponds to the classic form of the negative log marginal likelihood in a Bayesian linear regression context (see Bishop [5, Equation 3.86]).

Line (\star) introduces three terms that cancel out: $\frac{1}{2\sigma^2} \text{tr}(\Phi^T \Phi A^{-1}) + \frac{1}{2\sigma_\pi^2} \text{tr}(A^{-1}) - \frac{1}{2}d = 0$.

The latter equality follows from the trace operator properties and the definition of matrix A :

$$\begin{aligned} \frac{1}{2\sigma^2} \text{tr}(\Phi^T \Phi A^{-1}) + \frac{1}{2\sigma_\pi^2} \text{tr}(A^{-1}) &= \text{tr} \left(\frac{1}{2\sigma^2} \Phi^T \Phi A^{-1} + \frac{1}{2\sigma_\pi^2} A^{-1} \right) \\ &= \text{tr} \left(\frac{1}{2} A^{-1} \left(\frac{1}{\sigma^2} \Phi^T \Phi + \frac{1}{\sigma_\pi^2} \mathbf{I} \right) \right) \\ &= \text{tr} \left(\frac{1}{2} A^{-1} A \right) \\ &= \frac{1}{2} d. \end{aligned}$$

We show below that the expected loss $\mathbf{E}_{\mathbf{w} \sim \hat{\rho}^*} \hat{\mathcal{L}}_{X,Y}^{\text{nl}}(\mathbf{w})$ corresponds to the left part of Line (\star) . Note that a proof of equality $\mathbf{E}_{\mathbf{w} \sim \hat{\rho}} \mathbf{w}^T \Phi^T \Phi \mathbf{w} = \text{tr}(\Phi^T \Phi A^{-1}) + \hat{\mathbf{w}}^T \Phi^T \Phi \hat{\mathbf{w}}$ (Line \clubsuit below), known as the “expectation of the quadratic form”, can be found in Seber and Lee [49, Theorem 1.5].

$$\begin{aligned} n \mathbf{E}_{\mathbf{w} \sim \hat{\rho}} \hat{\mathcal{L}}_{X,Y}^{\text{nl}}(\mathbf{w}) &= \mathbf{E}_{\mathbf{w} \sim \hat{\rho}} \sum_{i=1}^n -\ln p(y_i | x_i, \mathbf{w}) \\ &= \mathbf{E}_{\mathbf{w} \sim \hat{\rho}} \left(\frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w} \cdot \phi(\mathbf{x}_i))^2 \right) \\ &= \frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \mathbf{E}_{\mathbf{w} \sim \hat{\rho}} \|\mathbf{y} - \Phi \mathbf{w}\|^2 \\ &= \frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \mathbf{E}_{\mathbf{w} \sim \hat{\rho}} (\|\mathbf{y}\|^2 - 2\mathbf{y} \Phi \mathbf{w} + \mathbf{w}^T \Phi^T \Phi \mathbf{w}) \\ &= \frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \left(\|\mathbf{y}\|^2 - 2\mathbf{y} \Phi \hat{\mathbf{w}} + \mathbf{E}_{\mathbf{w} \sim \hat{\rho}} \mathbf{w}^T \Phi^T \Phi \mathbf{w} \right) \\ &= \frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} (\|\mathbf{y}\|^2 - 2\mathbf{y} \Phi \hat{\mathbf{w}} + \text{tr}(\Phi^T \Phi A^{-1}) + \hat{\mathbf{w}}^T \Phi^T \Phi \hat{\mathbf{w}}) \quad (\clubsuit) \\ &= \frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{y} - \Phi \hat{\mathbf{w}}\|^2 + \frac{1}{2\sigma^2} \text{tr}(\Phi^T \Phi A^{-1}) \\ &= n \hat{\mathcal{L}}_{X,Y}^{\text{nl}}(\hat{\mathbf{w}}) + \frac{1}{2\sigma^2} \text{tr}(\Phi^T \Phi A^{-1}). \end{aligned}$$

Finally, the right part of Line (\star) is equal to the Kullback-Leibler divergence between the two multivariate normal distributions $\mathcal{N}(\hat{\mathbf{w}}, A^{-1})$ and $\mathcal{N}(\mathbf{0}, \sigma_\pi^2 \mathbf{I})$:

$$\begin{aligned} \text{KL}(\mathcal{N}(\hat{\mathbf{w}}, A^{-1}) \parallel \mathcal{N}(\mathbf{0}, \sigma_\pi^2 \mathbf{I})) &= \frac{1}{2} \left(\text{tr}((\sigma_\pi^2 \mathbf{I})^{-1} A^{-1}) + \frac{1}{\sigma_\pi^2} \|\hat{\mathbf{w}}\|^2 - d + \log \frac{|\sigma_\pi^2 \mathbf{I}|}{|A|} \right) \\ &= \frac{1}{2} \left(\frac{1}{\sigma_\pi^2} \text{tr}(A^{-1}) + \frac{1}{\sigma_\pi^2} \|\hat{\mathbf{w}}\|^2 - d + \log |A| + d \ln \sigma_\pi^2 \right). \end{aligned}$$

A.6 Linear Regression: PAC-Bayesian sub-gamma bound coefficients

We follow then exact same steps as in Section A.4, except that we replace the random variable v (giving the squared loss value) by a random variable v' giving the value of the loss

$$\ell_{\text{nl}}(\langle \mathbf{w}, \sigma \rangle, \mathbf{x}, y) = \frac{1}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2}(y - \mathbf{w} \cdot \mathbf{x})^2,$$

where \mathbf{w} , \mathbf{x} and y are generated as described in Section A.4. We aim to find the values of c and s such that the criterion of Equation (16) is fulfilled, *i.e.*,

$$\psi_{v'}(\lambda) = \ln \mathbf{E} e^{\lambda v'} \leq \frac{\lambda^2 s^2}{2(1-c\lambda)}, \quad \forall \lambda \in (0, \frac{1}{c}).$$

We obtain

$$\begin{aligned} \psi_{v'}(\lambda) &= \ln \mathbf{E}_{\mathbf{x}} \mathbf{E}_{y|\mathbf{x}} \mathbf{E}_{\mathbf{w}} \exp \left(\frac{\lambda}{2\sigma^2} \left[\mathbf{E}_{\mathbf{x}} \mathbf{E}_{y|\mathbf{x}} (y - \mathbf{w} \cdot \mathbf{x})^2 \right] - \lambda (y - \mathbf{w} \cdot \mathbf{x})^2 \right) \\ &\leq \ln \mathbf{E}_{\mathbf{w}} \exp \left(\frac{\lambda}{2\sigma^2} \mathbf{E}_{\mathbf{x}} \mathbf{E}_{y|\mathbf{x}} (y - \mathbf{w} \cdot \mathbf{x})^2 \right) \\ &\vdots \\ &= \frac{\frac{\lambda}{2\sigma^2} (\sigma_{\pi}^2 \sigma_{\mathbf{x}}^2 d + \sigma_{\mathbf{x}}^2 \|\mathbf{w}^*\|^2 + (1 - 2 \frac{\lambda}{2\sigma^2} \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2) \sigma_{\epsilon}^2)}{1 - 2 \frac{\lambda}{2\sigma^2} \sigma_{\mathbf{x}}^2 \sigma_{\pi}^2} \\ &= \frac{\lambda^2 s^2}{2(1 - \lambda c)}, \end{aligned} \tag{27}$$

with

$$\begin{aligned} c &= \frac{1}{2\sigma^2} \left[2\sigma_{\mathbf{x}}^2 \sigma_{\pi}^2 \right] = \frac{1}{\sigma^2} (\sigma_{\mathbf{x}}^2 \sigma_{\pi}^2), \\ s^2 &= \frac{1}{2\sigma^2} \left[\frac{2}{\lambda} \left[\sigma_{\mathbf{x}}^2 (\sigma_{\pi}^2 d + \|\mathbf{w}^*\|^2) + \sigma_{\epsilon}^2 (1 - \lambda c) \right] \right] = \frac{1}{\lambda \sigma^2} \left[\sigma_{\mathbf{x}}^2 (\sigma_{\pi}^2 d + \|\mathbf{w}^*\|^2) + \sigma_{\epsilon}^2 (1 - \lambda c) \right] \end{aligned}$$

Supplementary Material References

- [37] Luc Bégin, Pascal Germain, François Laviolette, and Jean-François Roy. PAC-Bayesian bounds based on the Rényi divergence. In *AISTATS*, pages 435–444, 2016.
- [38] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- [39] Peter Grünwald and John Langford. Suboptimal behavior of Bayes and MDL in classification under misspecification. *Machine Learning*, 66(2-3):119–149, 2007.
- [40] Peter Grünwald and Thijs van Ommen. Inconsistency of Bayesian Inference for Misspecified Linear Models, and a Proposal for Repairing It. *CoRR*, abs/1412.3730, 2014.
- [41] Peter D. Grünwald. *The Minimum Description Length Principle*. The MIT Press, 2007. ISBN 0262072815.
- [42] Sham M. Kakade and Andrew Y. Ng. Online bounds for Bayesian algorithms. In *NIPS*, pages 641–648, 2004.
- [43] Simon Lacoste-Julien, Ferenc Huszar, and Zoubin Ghahramani. Approximate inference for the loss-calibrated Bayesian. In *AISTATS*, pages 416–424, 2011.
- [44] John Langford and Matthias Seeger. Bounds for averaging classifiers. Technical report, Carnegie Mellon, Departement of Computer Science, 2001.
- [45] Ron Meir and Tong Zhang. Generalization error bounds for Bayesian mixture algorithms. *Journal of Machine Learning Research*, 4:839–860, 2003.
- [46] Andrew Y. Ng and Michael I. Jordan. On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes. In *NIPS*, pages 841–848. MIT Press, 2001.
- [47] Carl Rasmussen and Chris Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [48] Judith Rousseau. On the frequentist properties of Bayesian nonparametric methods. *Annual Review of Statistics and Its Application*, 3:211–231, 2016.
- [49] George A. F. Seber and Alan J. Lee. *Linear regression analysis*. John Wiley & Sons, 2012.