

# More PAC-Bayes bounds: From bounded losses, to losses with general tail behaviors, to anytime validity

**Borja Rodríguez-Gálvez**

BORJARG@KTH.SE

*Division of Information Science and Engineering (ISE)  
KTH Royal Institute of Technology  
Stockholm, Sweden*

**Ragnar Thobaben**

RAGNART@KTH.SE

*Division of Information Science and Engineering (ISE)  
KTH Royal Institute of Technology  
Stockholm, Sweden*

**Mikael Skoglund**

SKOGLUND@KTH.SE

*Division of Information Science and Engineering (ISE)  
KTH Royal Institute of Technology  
Stockholm, Sweden*

**Editor:** Ohad Shamir

## Abstract

In this paper, we present new high-probability PAC-Bayes bounds for different types of losses. Firstly, for losses with a bounded range, we recover a strengthened version of Catoni’s bound that holds uniformly for all parameter values. This leads to new fast-rate and mixed-rate bounds that are interpretable and tighter than previous bounds in the literature. In particular, the fast-rate bound is equivalent to the Seeger–Langford bound. Secondly, for losses with more general tail behaviors, we introduce two new parameter-free bounds: a PAC-Bayes Chernoff analogue when the loss’ cumulative generating function is bounded, and a bound when the loss’ second moment is bounded. These two bounds are obtained using a new technique based on a discretization of the space of possible events for the “in probability” parameter optimization problem. This technique is both simpler and more general than previous approaches optimizing over a grid on the parameters’ space. Finally, using a simple technique that is applicable to any existing bound, we extend all previous results to anytime-valid bounds.

**Keywords:** Generalization bounds, PAC-Bayes bounds, concentration inequalities, rate of convergence (fast, slow, mixed), tail behavior, parameter optimization.

## 1. Introduction

A learning algorithm  $\mathbb{A}$  is a (possibly randomized) mechanism that generates a hypothesis  $w \in \mathcal{W}$  of the solution of a certain problem given a sequence of  $n$  training data samples  $s := (z_1, \dots, z_n)$ , or *training set*. The performance of a hypothesis  $w$  on an instance  $z$  of the problem is described by a loss function  $\ell(w, z)$ . Hence, if the problem’s instances follow a distribution  $\mathbb{P}_Z$ , the goal is to produce a hypothesis  $w$  that attains a low *population risk*  $\mathcal{R}(w) := \mathbb{E}\ell(w, Z)$ , which is defined as the expected loss of the hypothesis  $w$  on samples  $Z$  drawn randomly from the problem’s distribution  $\mathbb{P}_Z$ .

Often, computing the population risk is not feasible. This is because, in general, the distribution  $\mathbb{P}_Z$  is unknown or intractable. However, having access to a training set  $s$ , a computable proxy for the population risk is the *empirical risk*  $\widehat{\mathcal{R}}(w, s) := \frac{1}{n} \sum_{i=1}^n \ell(w, z_i)$ , which is defined as the average loss of the hypothesis  $w$  on the samples from the training set  $s$ .

There are different attempts at characterizing the population risk based on the decomposition

$$\mathcal{R}(w) = \widehat{\mathcal{R}}(w, s) + \underbrace{\left( \mathcal{R}(w) - \widehat{\mathcal{R}}(w, s) \right)}_{\text{generalization error}}.$$

*Probably approximately correct (PAC)* theory gives bounds on the generalization error that hold with a probability larger than a certain threshold. Classically, these bounds depend only on the complexity of the hypothesis space  $\mathcal{W}$ , which is measured by, for example, the Vapnik–Cherovenkis (VC) dimension or the Rademacher complexity. See, for example, the book from Shalev-Shwartz and Ben-David (2014) for a pedagogical exposition of the topic.

In this work, we are concerned with *PAC-Bayesian bounds* (Shawe-Taylor et al., 1996; McAllester, 1998, 1999, 2003), which also consider the dependence of the hypothesis returned by the algorithm  $W = \mathbb{A}(S)$  on the training set  $S$ . These bounds are often of the following type: “for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \alpha_{\text{PAC-Bayes}}(S),”$$

where the probability is taken with respect to the sampling of the training set  $S \sim \mathbb{P}_S$  and  $\mathbb{E}^S$  denotes the conditional expectation operator with respect to the  $\sigma$ -algebra induced by  $S$ . The term  $\alpha_{\text{PAC-Bayes}}(S)$  describes the discrepancy between the population and empirical risks and is a random variable depending on  $S$ . Intuitively, this term (i) decreases with  $n$  as a better characterization of the risk is possible the more samples are available; (ii) increases with  $1/\beta$  as certainty comes with a price; and (iii) decreases as the hypothesis becomes less statistically dependent on the training set, as intuitively motivated by the fact that the empirical risk is always an unbiased estimate of the population risk in the extreme that the algorithm produces a fully independent hypothesis, that is  $\mathbb{E}\mathcal{R}(W) = \mathbb{E}\widehat{\mathcal{R}}(W, S)$ . In this paper, similarly to Hellström and Durisi (2020), we shall agree to the convention that the bounds are of *high probability* if the dependence on  $1/\beta$  is logarithmic, that is,  $\log 1/\beta$ . The review from Alquier (2021) offers an extensive introduction to PAC-Bayes theory.

McAllester (1998, 1999, 2003) showed the original PAC-Bayes bound considering bounded losses. The bound<sup>1</sup> states that for any prior  $\mathbb{Q}_W$  independent of  $S$  and every  $\beta \in (0, 1)$ , with probability no less than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \sqrt{\frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{2n}} \quad (1)$$

simultaneously for every Markov kernel  $\mathbb{P}_W^S$ , where  $\xi(n) \in [\sqrt{n}, 2 + \sqrt{2n}]$  (Maurer, 2004) and the dependency between the hypothesis and the dataset is measured by the relative

1. The bound written is the one obtained relaxing the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) via a lower bound on the binary relative entropy using Pinsker’s inequality. The term dependence  $\xi(n)$  with the number of samples  $n$  is the one established by Maurer (2004). See, for example, Section 2.2 of Tolstikhin and Seldin (2013).

entropy  $D(\mathbb{P}_W^S \|\mathbb{Q}_W)$  of the algorithm’s hypothesis kernel  $\mathbb{P}_W^S$ , or *posterior*, with respect to an arbitrary data-independent distribution on the hypothesis space  $\mathbb{Q}_W$ , or *prior*.<sup>2</sup> The dependency term  $D(\mathbb{P}_W^S \|\mathbb{Q})$  inside the square root plays the role of the complexity term of the classical PAC bounds, while the extra dependence  $\xi(n)$  on the number of samples comes from the concentration of the empirical risk around the population risk (see Maurer (2004) for the details). Finally, the term  $\log 1/\beta$  is the confidence penalty of being a high-probability bound. Sometimes, to simplify the discussion we will refer to this structure as the (normalized) *dependence-confidence* term and we define it as  $\mathfrak{C}_{n,\beta,S} := \frac{1}{n}(D(\mathbb{P}_W^S \|\mathbb{Q}_W) + \log 1/\beta)$ , which in the case of (1) corresponds to  $\mathfrak{C}_{2n,\beta/\xi(n),S}$ .

Many works on PAC-Bayes bounds have focused on two main tasks: (i) refining the bound to better characterize the population risk for bounded losses and (ii) extending this bound relaxing their assumptions or their setting.

In the first front, Langford and Seeger (2001); Seeger (2002) and Catoni (2003, 2007) developed more accurate bounds for estimating the population risk for bounded losses. However, either these bounds are not easily interpretable, minimizing them to find an appropriate posterior is hard, or they depend on an arbitrary parameter that needs to be selected *before* the draw of the data. To address these issues Tolstikhin and Seldin (2013), Thiemann et al. (2017), and Rivasplata et al. (2019) relaxed the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) to find more interpretable bounds where an approximate minimization to find a suitable posterior is possible. To contribute in this front:

**Contribution 1.** In Section 2, we show an alternative proof of a strengthened version of Catoni (2007)’s PAC-Bayes bound that holds uniformly for all values of the parameter  $\lambda$  (Theorem 6). We then build on this bound to show tighter fast-rate (Theorem 7) and mixed-rate (Theorem 9) bounds that are interpretable and help us to clarify the relationship between the population risk, the empirical risk, and the relative entropy of the algorithm’s posterior with respect to the prior. A mixed-rate bound is a bound with a mixture of a fast rate, and an amortized slow rate. The precise meaning becomes clear looking at Theorem 9. The fast-rate bound of Theorem 7 is of particular interest since it is equivalent to the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002). This reveals two significant insights: (i) that a linear combination of the empirical risk and the dependence-confidence term characterizes the bound and (ii) that the optimal posterior is a Gibbs distribution with a data-dependent “temperature”.

Wu and Seldin (2022) derived a “split-kl” inequality that competes with the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) for ternary losses and Jang et al. (2023) proved an even tighter bound via “coin-betting”. However, their bounds still neither are easily interpretable nor directly aid to the selection of an appropriate posterior. Moreover, there are other advances in this front when further quantities are considered. If the variance is known, Seldin et al. (2012, Theorem 8) and Wu et al. (2021, Theorem 9) introduced, respectively, PAC-Bayes analogues to Bernstein and Bennet inequalities. The

---

2. The range of  $\xi(n)$  is usually set to  $[\sqrt{n}, 2\sqrt{n}]$  for all  $n \geq 1$  as per the analysis of Maurer (2004) and empirical further analysis of Germain et al. (2015, Lemma 19). From Maurer (2004, Theorem 1), we can observe that the tighter  $2 + \sqrt{2n}$  is valid as an upper bound for all  $n \geq 2$  and the case where  $n = 1$  can be verified empirically using the bound from Germain et al. (2015, Lemma 19).

PAC-Bayes Bernstein inequality was later improved by further bounding the variance using an empirical estimate of that quantity (Tolstikhin and Seldin, 2013, Theorems 3 and 4). Finally, Mhammedi et al. (2019) derived a PAC-Bayes analogue to the “un-expected Bernstein inequality” where they use an empirical estimate of the second moment.

In the second front, Guedj and Pujol (2021) and Hellström and Durisi (2020) extended McAllester (2003)’s bound to subgaussian losses, resulting in the same rate as the original bound (1). However, the proof of these new bounds contains a small mistake. They derive intermediate PAC-Bayes bounds depending on a parameter  $\lambda$  that needs to be selected *before* the draw of the training data, and then they optimize this parameter without paying the necessary union bound price (Banerjee and Montúfar, 2021, Remark 14).<sup>3</sup> This leads to vacuous bounds as potentially an infinite number of parameters can be optimal for different data and it is a known standing problem in the PAC-Bayes literature (Alquier, 2021, Section 2.1.4). To address this issue:

**Contribution 2.** In Section 3, we devise a proof technique that allows us to bypass this optimization subtlety. We use this technique to extend McAllester (2003)’s bound to losses with more general tail behaviors. First, we derive a PAC-Bayes Chernoff analogue (Theorem 12) that specializes to the bounds of Guedj and Pujol (2021) and Hellström and Durisi (2020, 2021) for subgaussian losses. After that, we derive a parameter-free PAC-Bayes bound requiring only that the loss has a bounded second moment (Theorem 13). This last bound is of the nature of (Kuzborskij and Szepesvári, 2019) and is obtained by optimizing the parameter of Wang et al. (2015)’s bound on martingales in Appendix B.5. The proposed technique is simpler and more general than previous approaches that generate a grid on the parameters’ space, optimize the parameter over that grid, and pay the union bound price. Contrary to our technique, these approaches either can’t generate a parameter-free bound (Langford and Caruana, 2001; Catoni, 2003) or need to craft the grid in a case-to-case basis and need that an explicit solution of the optimal parameter exists (Seldin et al., 2012), which may not be the case (see Section 3.2.3). Therefore, the proposed technique is of independent interest for the development of future bounds “in probability”.

Other works also developed PAC-Bayes bounds with more general tail behaviors (Catoni, 2004; Alquier, 2006; Kuzborskij and Szepesvári, 2019; Haddouche and Guedj, 2023a; Alquier and Guedj, 2018; Holland, 2019; Haddouche et al., 2021; Chugg et al., 2023). However, most of these bounds either are not of high probability, or contain terms that often make the bounds non-decreasing with the number of samples  $n$ , or decrease at a slower rate than (1) when restricted to the bounded case, or also depend on parameters that need to be chosen *before* the draw of the training data.

Recently, some research has focused on developing PAC-Bayes bounds that hold *simultaneously* for all numbers of samples  $n$  (Chugg et al., 2023; Jang et al., 2023; Haddouche and Guedj, 2023a). These bounds are particularly useful for online learning algorithms that process data sequentially. These *anytime-valid* (or *time-uniform*) PAC-Bayes bounds are

---

3. Hellström and Durisi (2021) later corrected this issue using unique subgaussian properties (Wainwright, 2019, Theorem 2.6).

typically based on supermartingales and Ville (1939)’s extension of Markov’s inequality. To contribute in this end:

**Contribution 3.** In Section 4, we note that every PAC-Bayes bound can be extended to an anytime-valid one at a union bound cost (Theorem 16). For high-probability PAC-Bayes bounds, this cost is small.

Finally, note that while the relative entropy is widely used as the dependency measure in PAC-Bayes bounds due to its simplicity, interpretability from an information-theoretic perspective, and mathematical tractability, it has been shown that there are situations where an algorithm generalizes but this measure is large, making the bounds vacuous (Bassily et al., 2018; Livni and Moran, 2020; Haghifam et al., 2023; Nagarajan and Kolter, 2019; Nachum et al., 2023). The study of different dependency measures for PAC-Bayes bounds is outside the scope of this paper. Nonetheless, we refer the reader to other literature substituting the relative entropy as the dependency measure by different metrics like other  $f$ -divergences (Esposito et al., 2021; Ohnishi and Honorio, 2021; Kuzborskij et al., 2024), Rényi divergences (Bégin et al., 2016; Esposito et al., 2021; Hellström and Durisi, 2020), or integral probability metrics like the Wasserstein distance (Amit et al., 2022; Haddouche and Guedj, 2023b; Viallard et al., 2024a,b).

## 2. Specialized PAC-Bayes bounds for bounded losses

This section is separated into two parts. In Section 2.1, we review the state of the art of PAC-Bayes bounds for bounded losses. Then, in Section 2.2 we give an alternative proof of a strengthened version of Catoni (2007)’s parameterized bound that holds *simultaneously* for all values of the parameter. After that, we show that relaxing this strengthened bound (Theorem 6) yields fast-rate (Theorem 7 and Corollary 8) and mixed-rate (Theorem 9) bounds tighter than Thiemann et al. (2017)’s fast-rate and Tolstikhin and Seldin (2013)’s and Rivasplata et al. (2019)’s mixed-rate bounds.

### 2.1 A review of PAC-Bayes bounds for bounded losses

There are many important inequalities in the PAC-Bayes literature, especially for the case where the loss is bounded. These bounds are often presented for losses with a range in  $[0, 1]$ , which includes the interesting 0–1 loss for classification tasks. The Seeger–Langford (Langford and Seeger, 2001; Seeger, 2002) and Catoni (2007, Theorem 1.2.6)’s bounds are known to be (two of) the tightest bounds in this setting (cf. (Foong et al., 2021)). Both of them can be derived from Germain et al. (2009, Theorem 2.1)’s convex function bound. Below we state the extension from Bégin et al. (2014) that lifts the double absolute continuity requirement from the original statement noted by Haddouche et al. (2021).

**Theorem 1 (Bégin et al. (2016, Theorem 4))** *Consider a loss function  $\ell$  with bounded range  $[0, 1]$ , let  $\mathbb{Q}_W$  be any prior independent of  $S$ , and let  $W'$  be distributed according to  $\mathbb{Q}_W$ . Then, for every convex function  $f : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$  such that  $\mathbb{E}[\exp(nf(\widehat{\mathcal{R}}(W', s), \mathcal{R}(W')))] < \infty$  for all  $s \in \mathcal{Z}^n$ , and every  $\beta \in (0, 1)$ , with proba-*

bility no smaller than  $1 - \beta$

$$f(\mathbb{E}^S \widehat{\mathcal{R}}(W, S), \mathbb{E}^S \mathcal{R}(W)) \leq \frac{1}{n} \left[ D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{1}{\beta} + \log \mathbb{E} \left[ e^{nf(\widehat{\mathcal{R}}(W', S), \mathcal{R}(W'))} \right] \right]$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

This general bound is useful because an appropriate choice of the convex function  $f$  can be used to recover McAllester (2003)'s bound.<sup>4</sup> Similarly, choosing  $f(p, q) = d(p \| q) := D(\text{Ber}(p) \| \text{Ber}(q))$  combined with Maurer (2004)'s trick recovers the improved Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002), and choosing  $f(p, q) = -\log(1 - q(1 - e^{-\lambda/n})) - \lambda p/n$  recovers Catoni (2007, Theorem 1.2.6)'s bound.

**Theorem 2 (Improved Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002; Maurer, 2004))** Consider a loss function  $\ell$  with bounded range  $[0, 1]$  and let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$

$$d(\mathbb{E}^S \widehat{\mathcal{R}}(W, S) \| \mathbb{E}^S \mathcal{R}(W)) \leq \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} \quad (2)$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

**Theorem 3 (Catoni (2007, Theorem 1.2.6))** Consider a loss function  $\ell$  with bounded range  $[0, 1]$  and let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\lambda > 0$  and every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \frac{1}{1 - e^{-\frac{\lambda}{n}}} \left[ 1 - e^{-\frac{\lambda \mathbb{E}^S \widehat{\mathcal{R}}(W, S)}{n} - \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{1}{\beta}}{n}} \right]$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

The Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) is hindered by its lack of interpretability and the difficulty in minimizing it to find an appropriate posterior  $\mathbb{P}_W^S$ . This is due to the non-convexity of the bound with respect to the posterior  $\mathbb{P}_W^S$  (Thiemann et al., 2017) as well as the fact that it cannot be expressed explicitly as a function of the empirical risk  $\mathbb{E}^S \widehat{\mathcal{R}}(W, S)$  and the dependency term  $D(\mathbb{P}_W^S \| \mathbb{Q}_W)$  (Germain et al., 2009). On the other hand, while Catoni (2007)'s bound is minimized by the Gibbs posterior proportional to  $\mathbb{Q}_W(w) e^{-\lambda \widehat{\mathcal{R}}(w, S)}$ , it still lacks interpretability and depends on an arbitrary parameter  $\lambda$  that has to be chosen *before* the draw of the data.

To remedy these issues, several works relax the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) using lower bounds on the relative entropy (Tolstikhin and Seldin, 2013; Thiemann et al., 2017; Rivasplata et al., 2019). Since McAllester (2003)'s bound (1) is recovered with the standard Pinsker's inequality (Polyanskiy and Wu, 2023, Theorem 7.9), these works employ different relaxations of the stronger Marton (1996)'s bound, cf. (Seldin,

4. It is often mentioned that this is done choosing  $f(p, q) = 2(p - q)^2$ . However, technically, to use McAllester (2003)'s proof  $f(p, q) = (2^{n-1}/n)(p - q)^2$  should be used instead.

2023, Corollaries 2.19 and 2.20). Tolstikhin and Seldin (2013) use (Seldin, 2023, Corollary 2.20) and Thiemann et al. (2017) and Rivasplata et al. (2019) use (Seldin, 2023, Corollary 2.19). The latter bound results in an intractable PAC-Bayes bound; therefore Thiemann et al. (2017) relax it using the inequality  $\sqrt{xy} \leq \frac{1}{2}(\lambda x + y/\lambda)$  for all  $\lambda > 0$  to obtain a *fast-rate* bound, and Rivasplata et al. (2019) solve the resulting quadratic inequality for  $\sqrt{\mathbb{E}^S \mathcal{R}(W)}$  to obtain a *mixed-rate* bound.

**Theorem 4 (Thiemann et al. (2017, Theorem 3)’s fast-rate bound)** *Consider a loss function  $\ell$  with bounded range  $[0, 1]$  and let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \inf_{\lambda \in (0, 2)} \left\{ \frac{\mathbb{E}^S \widehat{\mathcal{R}}(W, S)}{1 - \frac{\lambda}{2}} + \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n\lambda(1 - \frac{\lambda}{2})} \right\}$$

*holds simultaneously for every posterior  $\mathbb{P}_W^S$ .*

**Theorem 5 (Rivasplata et al. (2019, Theorem 1)’s mixed-rate bound)** *Consider a loss function  $\ell$  with bounded range  $[0, 1]$  and let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$*

$$\begin{aligned} \mathbb{E}^S \mathcal{R}(W) &\leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} \\ &+ \sqrt{2\mathbb{E}^S \widehat{\mathcal{R}}(W, S) \cdot \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} + \left[ \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} \right]^2} \end{aligned}$$

*holds simultaneously for every posterior  $\mathbb{P}_W^S$ .*

Originally, Rivasplata et al. (2019) present their bound in a different form, but this form shows explicitly the combination of a *fast-rate* term and a *slow-rate* term. Moreover, this form makes it easy to see that the bound is tighter than (Tolstikhin and Seldin, 2013, Equation (3)) as their bound can be recovered using the inequality  $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$ .

## 2.2 From Seeger–Langford to an improved Catoni and new fast and mixed-rate bounds

As mentioned previously, both the Seeger–Langford (Langford and Seeger, 2001; Seeger, 2002) and Catoni (2007)’s bounds are known to be very tight (see Foong et al. (2021)) and are often used when a numerical certificate of the population risk is needed (Dziugaite and Roy, 2017; Pérez-Ortiz et al., 2021; Lotfi et al., 2022). Below, we show that a strengthened version of Catoni (2007)’s bound that holds *simultaneously* for all  $\lambda > 0$  can be obtained from the Seeger–Langford (Langford and Seeger, 2001; Seeger, 2002) bound at the small cost of  $\log \xi(n)$  in the dependence-confidence term. The idea behind the proof is to apply the Donsker and Varadhan (1975)’s lemma to  $d(\mathbb{E}^S \widehat{\mathcal{R}}(W, S) \| \mathbb{E}^S \mathcal{R}(W))$ . This was also observed by Germain et al. (2009, Proposition 2.1) and proved with different techniques than ours by (Catoni, 2015, Chapter 20) and Foong et al. (2021, Lemmata E1 and E2), although it was not stated explicitly as a PAC-Bayes bound.

**Theorem 6** Consider a loss function  $\ell$  with bounded range  $[0, 1]$  and let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \inf_{\lambda > 0} \left\{ \frac{1}{1 - e^{-\frac{\lambda}{n}}} \left[ 1 - e^{-\frac{\lambda \mathbb{E}^S \widehat{\mathcal{R}}(W, S)}{n} - \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n}} \right] \right\} \quad (3)$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

**Proof** Consider the Seeger–Langford bound from Theorem 2. Applying the Donsker and Varadhan (1975, Lemma 2.1)’s variational representation from Lemma 18 to the left hand side of (2) results in

$$\begin{aligned} d(\mathbb{E}^S \widehat{\mathcal{R}}(W, S) \| \mathbb{E}^S \mathcal{R}(W)) = \sup_{g_0, g_1 \in (-\infty, \infty)} \left\{ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) g_1 + (1 - \mathbb{E}^S \widehat{\mathcal{R}}(W, S)) g_0 \right. \\ \left. - \log \left[ \mathbb{E}^S \mathcal{R}(W) e^{g_1} + (1 - \mathbb{E}^S \mathcal{R}(W)) e^{g_0} \right] \right\}, \end{aligned}$$

where we defined  $g_0 := g(0)$  and  $g_1 := g(1)$ . Re-arranging the terms and plugging them into (2) states that with probability no smaller than  $1 - \beta$

$$\sup_{g_0, g_1 \in (-\infty, \infty)} \left\{ g_0 + \mathbb{E}^S \widehat{\mathcal{R}}(W, S) (g_1 - g_0) - \log \left[ e^{g_0} + \mathbb{E}^S \mathcal{R}(W) (e^{g_1} - e^{g_0}) \right] \right\} \leq \mathfrak{C}_{n, \beta / \xi(n), S}$$

holds *simultaneously* for every posterior  $\mathbb{P}_W^S$ . Note that, similarly to Thiemann et al. (2017)’s result from Theorem 4, the bound holds *simultaneously* for all values of the parameters  $g_0$  and  $g_1$ , and therefore these parameters can be chosen adaptively, that is, different values of  $g_0$  and  $g_1$  can be chosen for different realizations of the training set  $s$ . Therefore, with probability no smaller than  $1 - \beta$

$$-(g_0 - g_1) \mathbb{E}^S \widehat{\mathcal{R}}(W, S) - \log \left( 1 - (1 - e^{-(g_0 - g_1)}) \mathbb{E}^S \mathcal{R}(W) \right) \leq \mathfrak{C}_{n, \beta / \xi(n), S}$$

*simultaneously* for every posterior  $\mathbb{P}_W^S$  and all  $g_0$  and  $g_1$  in  $\mathbb{R}$ . Letting  $\lambda := n(g_0 - g_1)$  and re-arranging the terms in the equation it follows that, with probability no smaller than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \frac{1}{1 - e^{-\frac{\lambda}{n}}} \left[ 1 - e^{-\frac{\lambda \mathbb{E}^S \widehat{\mathcal{R}}(W, S)}{n} - \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n}} \right]$$

*simultaneously* for every posterior  $\mathbb{P}_W^S$  and all  $\lambda > 0$ . The restriction to  $\lambda > 0$  instead of  $\lambda \in \mathbb{R}$  comes from the fact that if  $\lambda < 0$ , then the resulting inequality is a lower bound instead of an upper bound.  $\blacksquare$

The bound in Theorem 6 is an explicit expression of the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) in terms of  $\mathbb{E}^S \widehat{\mathcal{R}}(W, S)$  and  $D(\mathbb{P}_W^S \| \mathbb{Q}_W)$ . Compared to Catoni (2007)’s Theorem 3, this bound holds *simultaneously* for all  $\lambda > 0$ , making it useful for finding numerical population risk certificates without the need to pay an extra price for the



parameter search. It also allows for an iterative procedure for obtaining a good posterior by updating the posterior  $\mathbb{P}_W^S$  and parameter  $\lambda$  alternately. We note that, contrary to the statment from the Seeger–Langford bound in Theorem 2, this statment tells us that the optimal posterior is given by the Gibbs distribution  $\mathbb{P}_W^S(w) \propto \mathbb{Q}_W(w) \cdot e^{-\lambda \widehat{\mathcal{R}}(w,S)}$ . However, finding the global optimum for the parameter  $\lambda$  is tedious, and the function is not convex in that parameter.

We may massage Theorem 6 to obtain a simpler, more interpretable fast-rate bound.

**Theorem 7** *Consider a loss function  $\ell$  with bounded range  $[0, 1]$  and let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \inf_{\substack{\gamma > 1 \\ c \in (0,1)}} \left\{ c\gamma \log \left( \frac{\gamma}{\gamma - 1} \right) \cdot \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + c\gamma \cdot \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} + \kappa(c)\gamma \right\} \quad (4)$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ , where  $\kappa(c) := 1 - c(1 - \ln c)$ .

**Proof** The bound follows by noting that the function  $1 - e^{-x}$  is a non-decreasing, concave, continuous function for all  $x > 0$  and therefore can be upper bounded by its envelope, that is,  $1 - e^{-x} = \inf_{a>0} \{e^{-a}x + 1 - e^{-a}(1 + a)\}$ . Using the envelope in (3), the changes of variable  $a := \log 1/c$  and  $\lambda := n \log \gamma/(\gamma-1)$ , and noting that  $c = e^{-a} \in (0, 1]$  and  $\gamma = (1 - e^{-\frac{\lambda}{n}})^{-1} > 1$  completes the proof.  $\blacksquare$

The parameter  $\gamma$  controls the influence of the empirical risk compared to the normalized dependence-confidence: if the empirical risk is large relative to the normalized dependence-confidence, then  $\gamma$  is larger and the normalized dependence-confidence coefficient increases, if instead the empirical risk is small or even close to interpolation, then  $\gamma$  is close to 1 and the empirical risk coefficient increases. In particular, for a fixed value of  $c$ , the optimal value of  $\gamma$  is

$$\begin{aligned} \gamma &= 1 + \left[ -1 - W \left( - \exp \left( -1 - \frac{c \cdot \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} + \kappa(c)}{c \cdot \mathbb{E}[\widehat{\mathcal{R}}(W, S)]} \right) \right) \right]^{-1} \\ &\approx 1 + \left[ \sqrt{2 \cdot \frac{c \cdot \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} + \kappa(c)}{c \cdot \mathbb{E}[\widehat{\mathcal{R}}(W, S)]}} + \frac{5}{6} \cdot \frac{c \cdot \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} + \kappa(c)}{c \cdot \mathbb{E}[\widehat{\mathcal{R}}(W, S)]} \right]^{-1}, \end{aligned}$$

where we considered the Lambert W function  $W$  and Chatzigeorgiou (2013)’s approximation of the  $-1$  branch.

The parameter  $c \in (0, 1]$  controls how much weight is given to the empirical risk and normalized dependence-confidence terms compared to a bias. For larger values of the empirical risk and the normalized dependence-confidence term, the value of  $c$  is small, decreasing their contribution to the bound and increasing the contribution of the bias  $\kappa(c) \in [0, 1)$ . If the empirical risk and the normalized dependence-confidence term are smaller, then the value

of  $c$  approaches 1, where the contribution of these two terms is only controlled by  $\gamma$  and the bias is 0. In fact, a weaker version of Theorem 7 can be obtained considering this small empirical risk and small normalized dependence-confidence regime by letting  $c = 1$ .

**Corollary 8** *Consider a loss function  $\ell$  with bounded range  $[0, 1]$  and let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \inf_{\gamma > 1} \left\{ \gamma \log \left( \frac{\gamma}{\gamma - 1} \right) \cdot \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \gamma \cdot \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} \right\} \quad (5)$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

Interestingly, this result can also be obtained using the variational representation of the relative entropy based on  $f$ -divergences (Polyanskiy and Wu, 2023, Example 7.5) as shown in Appendix A.1. That is, the Seeger–Langford bound (Theorem 2), the strengthened Catoni’s bound (Theorem 6), and this fast-rate bound (Theorem 7) are equally tight. This is important since it means that the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) can be exactly described with a linear combination of the empirical risk and the dependence-confidence term, where the coefficients of this combination and the bias vary depending on the data realization. This could have been hypothesized observing the derivatives of the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) from Reeb et al. (2018, Appendix A), and a proof is now available. Furthermore, the optimal posterior of this bound is given by the Gibbs distribution  $\mathbb{P}_W^S(w) \propto \mathbb{Q}_W(w) \cdot e^{-n \log \left( \frac{\gamma}{\gamma - 1} \right) \widehat{\mathcal{R}}(w, S)}$ , where the value of  $\gamma$  depends on the dataset realization  $s$ .

The bound improves upon Thiemann et al. (2017)’s Theorem 4 as it is tighter for all values of the empirical risk and the dependency measure (see Appendix A.2). For instance, the value  $\lambda = 1$  minimizes the multiplicative factor in the dependence-confidence term in Theorem 4. Letting  $\gamma = 2$  in Corollary 8 matches this factor and improves the multiplicative factor of the empirical risk from 2 to  $2 \log 2 \approx 1.38$ . Moreover, if we are in the *realizable setting* and  $\mathbb{E}^S \widehat{\mathcal{R}}(W, S) = 0$  (that is, we are using an empirical risk minimizer), then letting  $\gamma \rightarrow 1^+$  in this bound reveals that the fast rate can be achieved with multiplicative factor 1, clarifying that the dependence-confidence term completely characterizes the population risk in this regime. Note that this is not clear in Thiemann et al. (2017)’s nor Rivasplata et al. (2019)’s bounds, where the multiplicative factor is 2.

However, substituting the value of the optimal  $\gamma$  into (4) or (5) does not produce an interpretable bound. Nonetheless, the bound in Corollary 8 can be further relaxed to obtain a parameter-free mixed-rate bound that is tighter than Rivasplata et al. (2019)’s mixed-rate and Thiemann et al. (2017)’s fast-rate bounds (see Appendix A.2).

**Theorem 9 (mixed-rate bound)** *Consider a loss function  $\ell$  with bounded range  $[0, 1]$  and let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} + \sqrt{2 \mathbb{E}^S \widehat{\mathcal{R}}(W, S) \cdot \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n}} \quad (6)$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

**Proof** Using the inequality  $\log x \leq \frac{1}{2}(x - 1/x)$  for all  $x \geq 1$  in (5) establishes that with probability no smaller than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \inf_{\gamma > 1} \left\{ \frac{1}{2} \cdot \frac{2\gamma - 1}{\gamma - 1} \cdot \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \gamma \cdot \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{\xi(n)}{\beta}}{n} \right\}. \quad (7)$$

If  $\mathbb{E}^S \widehat{\mathcal{R}}(W, S) > 0$ , the optimal  $\gamma$ , which we recall can now be chosen adaptively as the bound holds simultaneously for all  $\gamma > 1$ , is  $\gamma = 1 + (\mathbb{E}^S \widehat{\mathcal{R}}(W, S) / 2\epsilon_{n, \beta/\xi(n), S})^{1/2}$ . Substituting this parameter yields (6). If  $\mathbb{E}^S \widehat{\mathcal{R}}(W, S) = 0$ , letting  $\gamma \rightarrow 1^+$  is optimal, which also agrees with the bound in (6).  $\blacksquare$

The mixed-rate bound presented in Theorem 9 provides a deeper insight into the relationship between the population risk, the empirical risk, and the dependence-confidence term. The bound grows linearly with both the empirical risk and the dependence-confidence term, with a correction term that reflects their interaction. Importantly, the bound is symmetric in these two terms, giving them equal importance. This may be beneficial for methods using PAC-Bayes bounds to optimize the posterior, such as PAC-Bayes with backprop (Rivasplata et al., 2019; Pérez-Ortiz et al., 2021), where using the bound from Thiemann et al. (2017) or Rivasplata et al. (2019) alone may cause the algorithm to disregard posteriors farther from the prior but that achieve lower population risk (see Appendix A.3).

### 3. PAC-Bayes bounds beyond bounded losses

This section is divided into two parts. In Section 3.1, we describe what a “more general tail behavior” means and review existing approaches to address this problem. In Section 3.2, we present a PAC-Bayes analogue to Chernoff’s inequality and a bound only requiring that the loss has a bounded second moment. Both bounds are obtained using a new technique for optimizing parameters that need to be selected *before* the draw of the data, but for which optimal solution depends on the data realization.

#### 3.1 What are losses with more general tail behaviors?

Probably, the most natural extension of a loss with a bounded range is a subgaussian loss (Wainwright, 2019, Chapter 2). The loss  $\ell(w, Z)$  is  $\sigma^2$ -subgaussian if it is concentrated around its mean with high probability. More precisely, it is *at least* as concentrated as if it were Gaussian with variance  $\sigma^2$ . That is, for all  $w \in \mathcal{W}$ , with probability no smaller than  $1 - \beta$  it holds that  $\mathbb{E}\ell(w, Z) \leq \ell(w, Z) + \sqrt{2\sigma^2 \log 1/\beta}$ . Notice that this definition of  $\sigma^2$ -subgaussian is equivalent to definitions expressed in terms of tail probabilities of the random variable  $\ell(w, Z)$ , or in terms of the moments of this random variable (see Wainwright (2019, Theorem 2.6)).

Hellström and Durisi (2020) and Guedj and Pujol (2021) obtained parameterized bounds for subgaussian losses. They also provided a parameter-free version of the bound optimizing the parameter. However, the optimization contained a small mistake, as the parameter needs to be selected *before* the draw of the data, and the value they chose depended on the data realization (Banerjee and Montúfar, 2021, Remark 14). Hellström and Durisi (2021) later

resolved this issue to obtain analogous PAC-Bayes bounds employing properties unique to subgaussian random variables (Wainwright, 2019, Theorem 2.6). Esposito et al. (2021) also derived PAC-Bayes bounds for this setting considering different dependency measures.

A step further beyond subgaussian losses are subexponential ones (Wainwright, 2019, Chapter 2). The loss  $\ell(w, Z)$  is subexponential if the probability that it is not concentrated around its mean is exponentially small. That is, for all  $w \in \mathcal{W}$ , with probability no smaller than  $1 - \beta$  it holds that  $\mathbb{E}\ell(w, Z) \leq \ell(w, Z) + 1/c_2 \log c_1/\beta$  for some  $c_1, c_2 > 0$ . Note that if a loss is subgaussian, then it is immediately subexponential by letting, for example,  $c_1 = e^{1/4}$  and  $c_2 = \sqrt{1/2\sigma^2}$ .

Catoni (2004) derived PAC-Bayes bounds for subexponential losses. However, these bounds are limited to the squared error loss in regression scenarios, where  $\mathcal{Z} = \mathcal{X} \times \mathbb{R}$  and the hypothesis  $w$  represents the parameters of a regressor  $\phi_w : \mathcal{X} \rightarrow \mathbb{R}$ . The analysis assumes that the regressor is finite, that is, that  $\|\phi_w\|_\infty < \infty$  for all  $w \in \mathcal{W}$ . Additionally, the derived bounds also rely on a parameter that must be chosen *before* the draw of the data.

These extensions can be generalized with the concept of *cumulant generating function* (CGF). The CGF of a random variable  $X$  is defined as  $\Lambda_X(\lambda) := \log \mathbb{E}[e^{\lambda(X - \mathbb{E}X)}]$  for all  $\lambda \in \mathbb{R}$ . If it exists, it completely characterizes the random variable's distribution, it is convex, continuously differentiable, and  $\Lambda_X(0) = \Lambda'_X(0) = 0$  (Zhang, 2006; Banerjee and Montúfar, 2021). We say that the CGF exists if it is bounded in  $(-b, b)$  for some  $b > 0$ .

**Definition 10 (Bounded CGF)** *A loss function  $\ell$  is of bounded CGF if for all  $w \in \mathcal{W}$ , there is a convex and continuously differentiable function  $\psi(\lambda)$  defined on  $[0, b)$  for some  $b \in \mathbb{R}_+$  such that  $\psi(0) = \psi'(0) = 0$  and  $\Lambda_{-\ell(w, Z)}(\lambda) \leq \psi(\lambda)$  for all  $\lambda \in [0, b)$ .*

Notice that we say that a loss has a bounded CGF if the CGF of  $-\ell(w, Z)$  is dominated by  $\psi$ . The reason is that we are interested in bounding from above the random variable  $(\mathcal{R}(w) - \ell(w, Z))$ , while the CGF of the loss  $\ell(w, Z)$  characterizes the tails of the random variable  $(\ell(w, Z) - \mathcal{R}(w))$ , where we recall that  $\mathcal{R}(w) = \mathbb{E}\ell(w, Z)$ .

Under this assumption, the *Cramér–Chernoff* method establishes that with probability no smaller than  $1 - \beta$  it holds that  $\mathbb{E}\ell(w, Z) \leq \ell(w, Z) + \psi_*^{-1}(\log 1/\beta)$  (Boucheron et al., 2013, Section 2.3), where  $\psi_*$  is the *convex conjugate* of the function  $\psi$  dominating the CGF, and  $\psi_*^{-1}$  is its inverse. More details about the convex conjugate and its inverse are given in Appendix B.1, but as illustrative examples, we have that for bounded losses  $\psi_*^{-1}(y) = \sqrt{y/2}$  and for  $\sigma^2$ -subgaussian losses  $\psi_*^{-1}(y) = \sqrt{2\sigma^2 y}$ .

Banerjee and Montúfar (2021) built on (Zhang, 2006) to prove a parameterized PAC-Bayes bound for losses with bounded CGF. The parameter must be chosen *before* drawing the data, similarly to the bounds in (Hellström and Durisi, 2020; Guedj and Pujol, 2021), which is a known standing issue in PAC-Bayes literature (Alquier, 2021). A slight extension is given below.

**Theorem 11 (Banerjee and Montúfar (2021, Theorem 6))** *Consider a loss function  $\ell$  with a bounded CGF in the sense of Definition 10. Let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\beta \in (0, 1)$  and every  $\lambda \in (0, b)$ , with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \frac{1}{\lambda} \left[ \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{1}{\beta}}{n} + \psi(\lambda) \right] \quad (8)$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

**Proof sketch** The proof follows similarly to the first part of Guedj and Pujol (2021, Proposition 3)’s proof with the addition of the ideas from Bu et al. (2020) and is given in Appendix B.2. ■

A further generalization of the tail behavior of the loss is to consider its moments. The *m-th moment* of the loss is defined as  $\mathbb{E}\ell(w, Z)^m$ . If the CGF exists (that is, it is bounded on an interval around zero), then all the moments are bounded. However, the reverse is not true: for instance, the Pareto distribution of the first kind with parameters  $a = 3$  and  $k = 1$  does not have a CGF but its variance is  $3/4$  (Norman L. Johnson, 1994, Chapter 20), and the lognormal distribution does not have a CGF but all its moments are finite (Norman L. Johnson, 1994, Chapter 14) (Asmussen et al., 2016). The smaller the moment assumed to be bounded, the weaker the assumption as if  $\mathbb{E}|\ell(w, Z)^m| < \infty$ , then  $\mathbb{E}|\ell(w, Z)^l| < \infty$  for all  $l \leq m$ .

Instead of directly bounding the population risk considering its tail behavior, Alquier (2006) proposed to bound it by studying a *truncated* version of the loss. That is, let  $\ell^-(w, z) := \min\{\ell(w, z), a\}$  and  $\ell^+(w, z) := \max\{\ell(w, z) - a, 0\}$  for some  $a \in \mathbb{R}$  such that  $\ell(w, z) \leq \ell^-(w, z) + \ell^+(w, z)$ . Then, one may find a PAC-Bayes bound for the population risk associated to this truncated loss  $\ell^-$  using the standard techniques from Section 2 and translate that into a PAC-Bayes bound for the real loss accounting for the tail probability  $\mathbb{E}^S \ell^+(w, Z)$ . Following this strategy, Alquier (2006) proposes bounds that still depend on a parameter that needs to be selected *before* the draw of the data. Moreover, similarly to Catoni (2004), for a regression problem with the square loss and a bounded regressor, Alquier (2006) shows the effect of the tail probability  $\mathbb{P}[\mathcal{A}_w^c]$  is not dominant even if the loss is subexponential.

Alquier and Guedj (2018) also developed PAC-Bayes bounds for losses with heavier tails that sometimes work for non i.i.d. data, although they are not of high probability and consider  $f$ -divergences as the dependency measure. Holland (2019) found PAC-Bayes bounds for losses with bounded second and third moments, but consider a different estimate than the empirical risk, and their bounds contain a term that may increase with the number of samples  $n$ . Finally, Kuzborskij and Szepesvári (2019) and Haddouche and Guedj (2023a) developed bounds for losses with a bounded second moment. The bound from Haddouche and Guedj (2023a) is anytime valid but depends on a parameter that needs to be chosen *before* the draw of the training data.

Haddouche et al. (2021) developed PAC-Bayes bounds under a different generalization, namely the hypothesis-dependent range (HYPE) condition, that is, that there is a function  $\kappa$  with positive range such that  $\sup_{z \in \mathcal{Z}} \ell(w, z) \leq \kappa(w)$  for all hypotheses  $w \in \mathcal{W}$ , but their bounds decrease at a slower rate than (1) when they are restricted to the bounded case. Finally, Chugg et al. (2023) also proved anytime-valid bounds for bounded CGFs and bounded moments, although their bounds contain parameters that need to be chosen *before* the draw of the training data with other technical conditions.

### 3.2 PAC-Bayes bounds for losses with bounded CGF or bounded second moment

McAllester (1998, 1999, 2003)'s PAC-Bayes bound (1) can be understood as a PAC-Bayes analogue to Hoeffding's inequality for bounded losses (Boucheron et al., 2013, Theorem 2.8): the two bounds are equivalent except for the dependency term between the hypothesis  $W$  and the training set  $S$ , namely  $D(\mathbb{P}_W^S \parallel \mathbb{Q}_W)$ . If we could optimize the parameter  $\lambda$  in Theorem 11, we would obtain a PAC-Bayes analogue to Chernoff's inequality for losses with a bounded CGF (Boucheron et al., 2013, Section 2.2). However, this is not possible since the optimal parameter depends on the data realization but needs to be selected *before* the draw of this data (Banerjee and Montúfar, 2021, Remark 14).

In the following theorem, we present a technique that allows us to bypass this subtlety for a small penalty of  $\log n$  or  $\log \log n$ . The idea is simple: separate the event space into a finite set of events where the optimization can be performed and then pay the union bound price. This can also be seen as optimizing over the set of parameters  $\lambda$  that will yield *almost optimal* bounds, and paying the union bound price for the cardinality of that set. In this case, the event space is separated using a quantization based on the relative entropy  $D(\mathbb{P}_W^S \parallel \mathbb{Q}_W)$  and noting that the event where  $D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) > n$  is not interesting as the resulting bound is non-decreasing with  $n$  given that event.

**Theorem 12 (PAC-Bayes Chernoff analogue)** *Consider a loss function  $\ell$  with a bounded CGF in the sense of Definition 10. Let  $\mathbb{Q}_W$  be any prior independent of  $S$  and define the event  $\mathcal{E} := \{D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) \leq n\}$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{1}_{\mathcal{E}} \cdot \left[ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) + \log \frac{en}{\beta}}{n} \right) \right] + \mathbb{1}_{\mathcal{E}^c} \cdot \text{ess sup } \mathbb{E}^S \mathcal{R}(W) \quad (9)$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ , where  $\mathbb{1}_{\mathcal{E}}$  is the indicator function defined as  $\mathbb{1}_{\mathcal{E}}(\omega) = 1$  if  $\omega \in \mathcal{E}$  and  $\mathbb{1}_{\mathcal{E}}(\omega) = 0$  otherwise.

**Proof** Let  $\mathcal{B}_\lambda$  be the complement of the event in (8) such that  $\mathbb{P}[\mathcal{B}_\lambda] < \beta$  and consider the sub-events  $\mathcal{E}_1 := \{D(\mathbb{P}_W^S \parallel \mathbb{Q}) \leq 1\}$  and  $\mathcal{E}_k := \{[D(\mathbb{P}_W^S \parallel \mathbb{Q})] = k\}$  for all  $k = 2, \dots, n$ , which form a covering of the event  $\mathcal{E}$ .<sup>5</sup> Furthermore, define  $\mathcal{K} := \{k \in \mathbb{N} : 1 \leq k \leq n \text{ and } \mathbb{P}[\mathcal{E}_k] > 0\}$ . For all  $k \in \mathcal{K}$ , given the event  $\mathcal{E}_k$ , with probability no more than  $\mathbb{P}[\mathcal{B}_\lambda | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \frac{1}{\lambda} \left[ \frac{k + \log \frac{1}{\beta}}{n} + \psi(\lambda) \right], \quad (10)$$

for all  $\lambda \in (0, b)$ . The right hand side of (10) can be minimized with respect to  $\lambda$  *independently of the training set  $S$* . Let  $\mathcal{B}_{\lambda_k}$  be the event resulting from this minimization and note that  $\mathbb{P}[\mathcal{B}_{\lambda_k}] \leq \beta$ . According to (Boucheron et al., 2013, Lemma 2.4), this ensures that with probability no more than  $\mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{k + \log \frac{1}{\beta}}{n} \right), \quad (11)$$

5. The notation  $\lceil x \rceil$  stands for the ceiling of  $x$ , that is, the nearest integer larger or equal to  $x$ .

where  $\psi_*$  is the convex conjugate of  $\psi$  and where  $\psi_*^{-1}$  is a non-decreasing concave function.<sup>6</sup> Given  $\mathcal{E}_k$ , since  $k < D(\mathbb{P}_W^S \|\mathbb{Q}_W) + 1$ , with probability no larger than  $\mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{D(\mathbb{P}_W^S \|\mathbb{Q}_W) + 1 + \log \frac{1}{\beta}}{n} \right).$$

Now, define  $\mathcal{B}'$  as the event stating that there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbf{1}_{\mathcal{E}} \cdot \left[ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{D(\mathbb{P}_W^S \|\mathbb{Q}_W) + \log \frac{\epsilon}{\beta}}{n} \right) \right] + \mathbf{1}_{\mathcal{E}^c} \cdot \text{ess sup } \mathbb{E}^S \mathcal{R}(W)$$

where  $\mathbb{P}[\mathcal{B}' | \mathcal{E}_k] \mathbb{P}[\mathcal{E}_k] \leq \mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k] \mathbb{P}[\mathcal{E}_k] \leq \mathbb{P}[\mathcal{B}_{\lambda_k}] \leq \beta$  for all  $k \in \mathcal{K}$  and where  $\mathbb{P}[\mathcal{B}' \cap \mathcal{E}^c] = 0$  by the definition of the essential supremum. Therefore, the probability of  $\mathcal{B}'$  is bounded as

$$\mathbb{P}[\mathcal{B}'] = \sum_{k \in \mathcal{K}} \mathbb{P}[\mathcal{B}' | \mathcal{E}_k] \mathbb{P}[\mathcal{E}_k] + \mathbb{P}[\mathcal{B}' \cap \mathcal{E}^c] < n\beta.$$

Finally, the substitution  $\beta \leftarrow \beta/n$  completes the proof.  $\blacksquare$

For subgaussian losses (and therefore for bounded ones), this recovers McAllester (2003)'s and Hellström and Durisi (2021)'s bound rates. For loss functions with heavier tails like *subgamma* and *subexponential*, the rates become a mixture of slow  $\sqrt{\mathfrak{C}_{n,\beta/en,S}}$  and fast  $\mathfrak{C}_{n,\beta/en,S}$  rates.<sup>7</sup> Please, see Corollary 24 in Appendix B.4 for an explicit formulation and derivation.

We can also employ this technique to obtain a parameter-free PAC-Bayes bound for losses with a bounded second moment optimizing the parameter in the results from Wang et al. (2015, Theorem 2.4) or Haddouche and Guedj (2023a, Theorem 2.1). The resulting bound is similar the one of Kuzborskij and Szepesvári (2019, Corollary 1), where the average sum of second moments plays the role of the subgaussian parameter in Theorem 12. Moreover, this bound essentially extends and improves upon the result from Alquier and Guedj (2018, Corollary 1) given the relationship between the relative entropy and the  $\chi^2$  divergence, that is, that  $D(\mathbb{P} \|\mathbb{Q}) \leq \log(1 + \chi^2(\mathbb{P} \|\mathbb{Q})) \leq \chi^2(\mathbb{P} \|\mathbb{Q})$  (Polyanskiy and Wu, 2023, Section 7.6). The proof and further details, including bounds for martingale sequences and non-i.i.d. data, are in Appendix B.5.

**Theorem 13 (Parameter-free bound for bounded second moment)** *Let  $\mathbb{Q}_W$  be any prior independent of  $S$  and define  $\xi'(n) := 2en(n+1)^2 \log(en)$  and the events  $\mathcal{E}_n := \{\sigma_n^2 D(\mathbb{P}_W^S \|\mathbb{Q}_W) \leq n\}$ , where  $\sigma_n^2 := \frac{1}{n} \sum_{i=1}^n \mathbb{E}^S[\ell(W, Z_i)^2 + 2\ell(W, Z')^2 + 1]$  for all  $n \in \mathbb{N}$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$*

$$\begin{aligned} \mathbb{E}^S \mathcal{R}(W) &\leq \\ &\mathbf{1}_{\mathcal{E}_n} \cdot \left[ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) \right] + \frac{2}{\sqrt{6}} \cdot \sqrt{\sigma_n^2 \left( \frac{D(\mathbb{P}_W^S \|\mathbb{Q}_W) + \log \frac{\xi'(n)}{\beta}}{n} \right)} + \mathbf{1}_{\mathcal{E}_n^c} \cdot \text{ess sup } \mathbb{E}^S \mathcal{R}(W). \end{aligned} \tag{12}$$

6. The infimum is not always attained with a particular value  $\lambda_k$  of  $\lambda$ . The details are given in Appendix B.3.

7. Subgamma random variables are nearly subgaussian, but not quite (Boucheron et al., 2013, Section 2.4).

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

Note that in the bound from Theorem 13, the terms  $\mathbb{E}^S \ell(W, Z_i)^2$  are fully empirical and the term  $\mathbb{E}^S \ell(W, Z')^2$  accounts for the assumption that the second moment of the loss is bounded. Theorem 13 is more general than Theorem 12, as only the knowledge of one moment is required instead of the knowledge of a function dominating the CGF, which signifies information of all the moments. However, the resulting bound is not always better. For instance, for  $\sigma^2$ -subgaussian,  $(\sigma^2, c)$ -subgamma, or  $(\sigma^2, c)$ -subexponential losses, the parameter  $\sigma^2$  that appears in Corollary 24 to Theorem 12 (cf. Appendix B.4) is a proxy for the variance  $\mathbb{V}^S \ell(W, Z') = \mathbb{E}^S (\ell(W, Z') - \mathbb{E}^S \mathcal{R}(W))^2$  or the *central* second moment of the loss, while  $\mathbb{E}^S \ell(W, Z')^2$  is its *raw* second moment, which can be much larger than the variance because  $\mathbb{E}^S \ell(W, Z')^2 = \mathbb{V}^S \ell(W, Z') + (\mathbb{E}^S \mathcal{R}(W))^2$ .

### 3.2.1 SMALLER UNION BOUND COST

Similarly to Langford and Seeger (2001) and Catoni (2003), we can pay a multiplicative cost of  $e$  to the relative entropy to reduce the union bound cost to  $\log^{(2+\log n)}/\beta$ . For example, for Theorem 12 the idea is to follow its proof with the events  $\mathcal{E}_k := \{\mathbb{D}(\mathbb{P}_W^S \| \mathbb{Q}_W) \in (e^{k-1}, e^k]\}$  and note that  $e^k < e\mathbb{D}(\mathbb{P}_W^S \| \mathbb{Q}_W)$  given  $\mathcal{E}_k$ . As mentioned by Maurer (2004), however, these bounds are only useful when the dependency measure  $\mathbb{D}(\mathbb{P}_W^S \| \mathbb{Q}_W)$  grows slower than logarithmically. This procedure is detailed in Appendix B.6.

### 3.2.2 DIFFERENT OR ABSENCE OF UNINTERESTING EVENTS

Theorems 12 and 13 consider the events  $\{\mathbb{D}(\mathbb{P}_W^S \| \mathbb{Q}_W) \leq n\}$  and  $\{\sigma_n^2 \mathbb{D}(\mathbb{P}_W^S \| \mathbb{Q}_W) \leq n\}$  since the complementary events are uninteresting as the bounds become  $\Omega(1)$ . However, if one is interested in a different event such as  $\{\mathbb{D}(\mathbb{P}_W^S \| \mathbb{Q}_W) \leq k_{\max}\}$  or  $\{\sigma_n^2 \mathbb{D}(\mathbb{P}_W^S \| \mathbb{Q}_W) \leq k_{\max}\}$ , then the proofs may be replicated. The resulting bounds are equal to (9) and (12), where the factor inside the logarithm will be  $e^{k_{\max}}/\beta$  and  $\xi^{(k_{\max})}/\beta$  respectively. Some examples would be to choose  $k_{\max} = \lceil \log(dn) \rceil$  for parametric models such as the sparse single-index (Alquier and Biau, 2013) and sparse additive (Guedj and Alquier, 2013) models, where  $d$  is the dimension of the input data, or to choose  $k_{\max} = \lceil \log(dpn) \rceil$  for the noisy  $d \times p$  matrix completion problem (Mai and Alquier, 2015). Other examples are given in Appendices B.4 and B.8.

Imagine that one is interested in a bound like those presented in Theorems 12 and 13 and does not consider any event to be uninteresting. This could happen in some regression application where even if  $\mathbb{D}(\mathbb{P}_W^S \| \mathbb{Q}_W) \geq n$  and the bound is in  $\Omega(1)$  the particular value of the bound is necessary. In this case, working in the events' space is still beneficial. The idea is almost the same as before: separate the events' space into a countable set of events where the optimization can be performed and pay the union bound price. The main difference is that each of these events  $\mathcal{E}_k$  will be defined with a different value of  $\beta_k$  so that price of the union bound is finite  $\sum_{k=1}^{\infty} \beta_k < \infty$ . For instance, applying this approach to Theorem 11 results in the following theorem, whose proof is in Appendix B.7.

**Theorem 14** *Consider a loss function  $\ell$  with a bounded CGF in the sense of Definition 10. Let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller*



than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{e\pi^2 (D(\mathbb{P}_W^S \| \mathbb{Q}_W) + 1)^2}{6\beta}}{n} \right)$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

Since  $x + \log e\pi^2(x+1)^2/6\beta$  is a non-decreasing, concave, continuous function for all  $x > 0$ , it can be upper bounded by its envelope. That is,  $x + \log e\pi^2(x+1)^2/6\beta \leq \inf_{a>0} \left(\frac{a+3}{a+1}\right)x + \log e\pi^2(a+1)^2/6\beta - \frac{2a}{a+1}$ . Taking  $a = 19$  leads to the following corollary.

**Corollary 15** *Consider a loss function  $\ell$  with a bounded CGF in the sense of Definition 10. Let  $\mathbb{Q}_W$  be any prior independent of  $S$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{1.1D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{10e\pi^2}{\beta}}{n} \right).$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

### 3.2.3 RELATED WORK

A related, but different technique to deal with these optimization problems is given by Langford and Caruana (2001) and Catoni (2003) to solve the bounded losses analogue of (11). They consider the optimization of  $\lambda$  over a geometric grid  $\mathcal{A} = \{e^k : k \in \mathbb{N}\} \cap [1, n]$  at the smaller union bound cost of  $\log^{(1+\log n)}/\beta$  at the price of a multiplicative constant of  $e$ . Using rounding arguments similar to those in the proofs of Theorems 12 and 13, this translates into being able to optimize the parameter  $\lambda$  in the region  $[1, n]$ . This technique generalizes to other countable families  $\mathcal{A}$  with a union bound cost of  $\log |\mathcal{A}|$  (Alquier, 2021, Section 2.1.4). The downfall of this approach compared to the one presented here is that the optimal parameter  $\lambda^*$  is still dependent on the data drawn  $S$ , the probability parameter  $\beta$ , and the tail behavior captured either by  $\psi_*^{-1}$  or  $\sigma_n^2$ . It is hence uncertain if the optimal parameter will lie within the set  $\mathcal{A}$  in general, making a parameter-free expression for the bound impossible.

An extension of this technique is given by Seldin et al. (2012). The idea is to construct a countably infinite grid  $\mathcal{A}$  over the parameters' space and then choose a parameter  $\lambda$  from that grid. Then, they can give a closed-form solution by studying how far is the bound resulting from plugging the selected parameter from the grid and the optimal parameter. Their technique has been used for a bounded range and bounded variance setting in (Seldin et al., 2012) and for a bounded empirical variance in (Tolstikhin and Seldin, 2013).

The main difference between these approaches and ours is that they design a grid  $\mathcal{A}$  over the parameters' space and optimize the parameter  $\lambda$  in that grid. Then, the tightness of the resulting bound depends on how well that grid was crafted. This grid  $\mathcal{A}$  needs to be designed in a case-to-case basis and it can be cumbersome, see for example (Tolstikhin and Seldin, 2013, Appendix A). Moreover, to design the said grid one requires an explicit expression for the optimal parameter. This may not be available in cases such as in Theorem 12, where we

only know that (11) is the result of the optimization in (10). On the other hand, we consider a grid over the events’ space and find the *best* parameter for each cell (sub-event) in that grid. This gives three main advantages with respect to the previous techniques. First, the grid is the same for any situation, making the technique easier to employ (see Theorems 12 to 14, Theorems 27 and 28, and Corollary 24). For instance, it would be trivial to recover a result similar to the PAC-Bayes Bernstein analog of Seldin et al. (2012, Theorem 8) optimizing the parameter in (Seldin et al., 2012, Theorem 7) with our approach. Second, to apply the technique we do not need to know the explicit form of the optimal parameter, which may not exist like in Theorem 12, we only need that the optimization is possible. Third, if the grid is made with respect to a random variable  $X$ , the resulting bound will be tight except from a logarithmic term and an offset changing  $X$  by  $X + 1$ . Therefore, discretizing the events’ space is essentially equivalent to crafting a subset  $\mathcal{A}'$  of the parameters’ space (not necessarily with a grid structure) with the *optimal* parameters for each region without the need to design this subset  $\mathcal{A}'$  in a case-to-case basis.

Another possibility to deal with  $\lambda$  is to integrate it with respect to an analytically integrable probability density with mass concentrated in its maximum. This is the method employed by Kuzborskij and Szepesvári (2019) and is known as *the method of mixtures* (de la Peña et al., 2007, Section 2.3). Unfortunately, this method requires the existence of a canonical pair: two random variables  $X$  and  $Y$  satisfying that  $\mathbb{E} \exp(\lambda X + \lambda^2 Y^2/2) \leq 1$  for all  $\lambda$  in the domain of optimization (de la Peña et al., 2007, Equation (2.2)). This requirement may not necessarily hold in general settings like Theorem 12. Moreover, often this method results in the introduction of a new parameter associated with the density used for integration, for example, the variance of a Gaussian as in (Kuzborskij and Szepesvári, 2019). Therefore, our proposed approach is still more general while resulting in essentially the same bound when restricted to the case where the method of mixtures can be employed.<sup>8</sup>

Finally, Kakade et al. (2008, Corollary 8) employed a similar technique to ours to prove a PAC-Bayes bound for bounded losses similar to McAllester (1998, 1999, 2003)’s Equation (1). However, they did not employ the technique to optimize a parameter. Instead, they found a bound in terms of a threshold  $a$  that held for every posterior  $\mathbb{P}_W^S$  such that  $D(\mathbb{P}_W^S \parallel \mathbb{Q}) \leq a$ . Then, they discretized the set of all posteriors into the sub-classes  $\mathcal{P}_k := \{\mathbb{P}_W^S : 2^{k+1} < D(\mathbb{P}_W^S \parallel \mathbb{Q}) \leq 2^{k+2}\}$  and applied the union bound to find a uniform result. This technique is usually known as the *peeling device*, *stratification*, or *slicing* in the probability theory and bandits communities (Boucheron et al., 2013, Section 13.7) (Lattimore and Szepesvári, 2020, Section 9.1). The similarity with our proof of Theorems 12, 14 and 28 is clear by looking at our design of the events’ discretization and their posterior’s sub-classes. However, the nature of the two approaches is different: they have a natural constraint, and they discretize the posterior class space and apply the union bound to circumvent that; while we have a parameter which optimal value is data-dependent, we discretize the events’ space to find the optimal parameter in a data-independent way, and then we apply the union bound. Moreover, this technique is more general, as one can design the sub-events to include basically any random object that depends on the data as showcased in Theorems 13 and 27. Nonetheless, one could consider our technique to be essentially equivalent to the peeling device, since both techniques have the same idea and intention behind them.

---

8. The final bounds are not directly comparable due to differences in the logarithmic terms, but both are of the same order.

### 3.2.4 IMPLICATIONS TO THE DESIGN OF POSTERIOR DISTRIBUTIONS

We will focus the discussion about which are the implications of having a parameter-free bound with more general assumptions with respect to the design of posterior distributions to Theorem 12. The discussion extends to Theorems 13 and 14 and other situations analogously.

The first consideration is that the parameter-free bound in (9) can always be transformed back into a parametric bound that holds *simultaneously* for all parameters. In the case of Theorem 12, employing Lemma 22 we have that with probability no smaller than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{1}_{\mathcal{E}} \cdot \left[ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \inf_{\lambda \in (0, b)} \left\{ \frac{D(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log \frac{en}{\beta} + \frac{\psi(\lambda)}{\lambda}}{\lambda n} \right\} \right] + \mathbb{1}_{\mathcal{E}^c} \cdot \text{ess sup } \mathbb{E}^S \mathcal{R}(W)$$

holds *simultaneously* for every posterior  $\mathbb{P}_W^S$ . This relaxation to a familiar structure tells us that the optimal posterior is the Gibbs distribution  $\mathbb{P}_W^S(w) \propto \mathbb{Q}_W(w) \cdot e^{\lambda n \widehat{\mathcal{R}}(w, S)}$ , where the value of  $\lambda$  can now be chosen *adaptively* for each dataset realization  $s$ .

The second consideration is if we are using some numerical estimation of the posterior using neural networks as with the PAC-Bayes with backprop (Rivasplata et al., 2019; Pérez-Ortiz et al., 2021) or other similar frameworks (Dziugaite and Roy, 2017; Lotfi et al., 2022). Then, the posterior can be readily estimated as long as the inverse of the convex conjugate  $\psi_*^{-1}$  is a differentiable function.

## 4. Anytime-valid PAC-Bayes bounds

Recently, some works (Chugg et al., 2023; Jang et al., 2023; Haddouche and Guedj, 2023a) have focused on *anytime valid* (or *time-uniform*) PAC-Bayes bounds, that is, bounds that hold *simultaneously* for all number of samples  $n$ . Often, their goal is to provide guarantees at every step for online algorithms that are sequential in nature. These bounds are usually rooted in the usage of supermartingales and Ville (1939)’s extension of Markov’s inequality.

Every standard PAC-Bayes bound can be extended to an anytime-valid bound at a union bound cost, even if it does not have a suitable supermartingale structure. For high-probability PAC-Bayes bounds like Theorems 12 and 13, this extension comes at the small cost of adding  $2 \log \pi n / \sqrt{6}$  to  $\log en / \beta$  or to  $\log \xi^{(n)} / \beta$  respectively. This “folklore” result is formalized below for general probabilistic bounds. Similar uses of the union bound in other settings appear in (Darling and Robbins, 1967; Robbins and Siegmund, 1968; Kaufmann et al., 2016; Howard et al., 2021).

**Theorem 16 (From standard to anytime-valid bounds)** *Consider the probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  and let  $(\mathcal{E}_n)_{n=1}^\infty$  be a sequence of event functions such that  $\mathcal{E}_n : (0, 1) \rightarrow \mathcal{F}$ . If  $\mathbb{P}[\mathcal{E}_n(\beta)] \geq 1 - \beta$  for all  $\beta \in (0, 1)$  and all  $n \geq 1$ , then  $\mathbb{P}[\cap_{n=1}^\infty \mathcal{E}_n(6\beta/\pi^2 n^2)] \geq 1 - \beta$  for all  $\beta \in (0, 1)$ .*

**Proof** We prove the equivalent statement: “for every  $\beta \in (0, 1)$ , if  $\mathbb{P}[\mathcal{E}_n^c(\beta)] < \beta$  for all  $n \geq 1$ , then  $\mathbb{P}[\cup_{n=1}^\infty \mathcal{E}_n^c(6\beta/\pi^2 n^2)] < \beta$ ”. By the union bound, it follows that  $\mathbb{P}[\cup_{n=1}^\infty \mathcal{E}_n^c(\beta_n)] < \sum_{n=1}^\infty \beta_n$ . Let  $\beta_n = \beta/n^2$ , then  $\mathbb{P}[\cup_{n=1}^\infty \mathcal{E}_n^c(\beta/n^2)] < \pi^2 \beta / 6$ . The substitution  $\beta \leftarrow 6\beta/\pi^2$  completes the proof.  $\blacksquare$

There are better choices of  $\beta_n$  such as  $\beta_n = \beta/n \log^2(6n)$  (Kaufmann et al., 2016), but all result in essentially the same cost  $\mathcal{O}(\log n)$  for high-probability PAC-Bayes bounds. The main takeaway from this result is that the anytime-valid bounds obtained via supermartingales and Ville (1939)’s inequality only contribute in shaving-off a log factor for PAC-Bayes high-probability bounds. Hence, their main advantage is in describing online learning situations where the subsequent samples are dependent to each other, which is not inherently captured by statements like Theorem 16.

**Remark 17** *Theorems 6, 7 and 9 and Corollary 8 follow verbatim as an anytime-valid bound substituting  $\log \xi^{(n)}/\beta$  by  $\log \sqrt{\pi^{(n+1)}}/\beta$  without needing Theorem 16. The reason is that these results are derived from the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002; Maurer, 2004), which is extended to an anytime-valid bound at this cost in (Jang et al., 2023).*

## 5. Conclusion

In this paper, we present new high-probability PAC-Bayes bounds. For bounded losses, the strengthened version of Catoni’s bound (Theorem 6) provides tighter fast and mixed-rate bounds (Theorems 7 and 9 and Corollary 8). Moreover, the fast-rate bound is equivalent to the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002), helping us to better understand the behavior of this bound and its optimal posterior. Namely, this reveals that the bound is completely characterized by a linear combination of the empirical risk and the dependence-confidence term, and that the optimal posterior of the fast-rate bound is a Gibbs distribution with a data-dependent “temperature”. For more general losses, we introduce two parameter-free bounds using a new technique to optimize parameters in probabilistic bounds: one for losses with bounded CGF and another one for losses with bounded second moment (Theorems 12 and 13). We also extend all our results to anytime-valid bounds with a technique that can be applied to any existing bound (Theorem 16).

PAC-Bayes bounds have been proven useful to provide both numerical population risk certificates as well as to understand the sufficient conditions for a problem to be learned (Ambroladze et al., 2006; Ralaivola et al., 2010; Higgs and Shawe-Taylor, 2010; Seldin and Tishby, 2010; Alquier and Biau, 2013; Guedj and Alquier, 2013; Mai and Alquier, 2015; Appert and Catoni, 2021; Nozawa and Sato, 2019; Nozawa et al., 2020; Chérif-Abdellatif et al., 2022). The interpretability of the bounds in Section 2 and the wider applicability of those in Section 3, along with their potential extension from Section 4, can contribute to extend this understanding. However, it is known that there are situations where an algorithm generalizes but the dependency measure (relative entropy) in the PAC-Bayes bounds is large, yielding them vacuous (Bassily et al., 2018; Livni and Moran, 2020; Haghifam et al., 2023). Nonetheless, some approaches recently managed to use PAC-Bayes bounds to obtain non-vacuous bounds for neural networks (Dziugaite and Roy, 2017, 2018; Rivasplata et al., 2019; Pérez-Ortiz et al., 2021; Zhou et al., 2019; Lotfi et al., 2022). The bounds from Section 2 can contribute in this front via methods like PAC-Bayes with backprop (Rivasplata et al., 2019; Pérez-Ortiz et al., 2021) as they are differentiable and tighter than previous bounds of this kind.

Technically, the procedure employed in Section 3, focusing on discretizing the event space instead of the parameter space, can be of independent interest and useful for developing

theory elsewhere. Similarly, Theorem 16 is presented in a general way so it can be seemingly used in other contexts beyond PAC-Bayesian theory.

## Acknowledgments

First of all, we would like to thank the reviewers and the editor for their comments and suggestions, which helped us improve the paper.

We are also grateful to Gergely Neu and our fruitful discussions that lead to a cleaner exposition of Theorem 12 using indicator functions; to Pierre Alquier for the suggestion of improving Section 3.2.2 considering a cut-off appropriate for PAC-Bayes bounds for parametric problems and for pointing us to further PAC-Bayes literature; and to Omar Rivasplata and María Pérez-Ortiz for their help dealing with PAC-Bayes with backprop.

This work was funded, in part, by the Swedish research council under contracts 2019-03606 (Borja Rodríguez-Gálvez and Mikael Skoglund) and 2021-05266 (Ragnar Thobaben).

## Appendix A. Details of Section 2

This section of the appendix is devoted to providing alternative proofs and supplementary context and examples to the results from Section 2.

First of all, to aid with the proofs of Theorem 6 and Theorem 7, we state the Donsker–Varadhan (Donsker and Varadhan, 1975) and  $f$ -divergence-based variational representations of the relative entropy (Polyanskiy and Wu, 2023, Example 7.5). For the next two lemmata,  $\mathcal{X}$  denotes a measurable space.

### Lemma 18 (Donsker–Varadhan variational representation of the relative entropy)

Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two probability measures on  $\mathcal{X}$  and  $X \sim \mathbb{P}$  and  $Y \sim \mathbb{Q}$  be two random variables. Further let  $\mathcal{G}$  be the set of functions  $g : \mathcal{X} \rightarrow \mathbb{R}$  such that  $\mathbb{E}[e^{g(Y)}] < \infty$ . Then,

$$D(\mathbb{P} \parallel \mathbb{Q}) = \sup_{g \in \mathcal{G}} \left\{ \mathbb{E}[g(X)] - \log \mathbb{E}[e^{g(Y)}] \right\}.$$

### Lemma 19 ( $f$ -divergence-based variational representation of the relative entropy)

Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two probability measures on  $\mathcal{X}$  and  $X \sim \mathbb{P}$  and  $Y \sim \mathbb{Q}$  be two random variables. Further let  $\mathcal{G}$  be the set of functions  $g : \mathcal{X} \rightarrow \mathbb{R}$  such that  $\mathbb{E}[e^{g(Y)}] < \infty$ . Then,

$$D(\mathbb{P} \parallel \mathbb{Q}) = \sup_{g \in \mathcal{G}} \left\{ \mathbb{E}[g(X)] - \mathbb{E}[e^{g(Y)}] + 1 \right\}.$$

### A.1 Alternative proof of Theorem 7

As mentioned in Section 2.2, Theorem 7 can be recovered directly from the Seeger–Langford (Langford and Seeger, 2001; Seeger, 2002) bound from Theorem 2 using the variational representation of the relative entropy based on  $f$ -divergences (Polyanskiy and Wu, 2023, Example 7.5).

**Alternative proof of Theorem 7** Similarly to the proof of Theorem 6, the proof starts considering the Seeger–Langford bound from Theorem 2. Now, applying the variational

representation of the relative entropy based on  $f$ -divergences from Lemma 19 to the right hand side of (2) results in

$$d(\mathbb{E}^S \widehat{\mathcal{R}}(W, S) \| \mathcal{R}(W)) = \sup_{g_0, g_1 \in (-\infty, \infty)} \left\{ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) g_1 + (1 - \mathbb{E}^S \widehat{\mathcal{R}}(W, S)) g_0 - \mathbb{E}^S \mathcal{R}(W) e^{g_1} + (1 - \mathbb{E}^S \mathcal{R}(W)) e^{g_0} + 1 \right\},$$

where we defined  $g_0 := g(0)$  and  $g_1 := g(1)$ . Re-arranging the terms and plugging them into (2) states that with probability no smaller than  $1 - \beta$

$$\sup_{g_0, g_1 \in (-\infty, \infty)} \left\{ 1 + g_0 + \mathbb{E}^S \widehat{\mathcal{R}}(W, S) (g_1 - g_0) - e^{g_0} - \mathbb{E}^S \mathcal{R}(W) (e^{g_1} - e^{g_0}) \right\} \leq \mathfrak{C}_{n, \beta/\xi(n), S}.$$

where we recall that  $\mathfrak{C}_{n, \beta, S} := \frac{1}{n} (\mathbb{D}(\mathbb{P}_W^S \| \mathbb{Q}_W) + \log 1/\beta)$ . Again, similarly to Thiemann et al. (2017)'s result from Theorem 4, the bound holds uniformly over all realizations of the training set, and thus the parameters  $g_0$  and  $g_1$  can be chosen adaptively. For this inequality to be relevant to us, we require that  $g_0 \geq g_1$ , as otherwise we would obtain a lower bound instead of an upper bound. To simplify the equations, let  $\gamma := e^{g_0}/(e^{g_0} - e^{g_1}) \geq 1$ , which implies that  $g_0 - g_1 = \log(\gamma/(\gamma - 1))$  and therefore with probability no smaller than  $1 - \beta$

$$1 + g_0 - e^{g_0} - \log\left(\frac{\gamma}{\gamma - 1}\right) \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \gamma^{-1} e^{g_0} \mathbb{E}^S \mathcal{R}(W)$$

*simultaneously* for every posterior  $\mathbb{P}_W^S$  and all  $g_0$  and  $g_1$  in  $\mathbb{R}$  such that  $g_0 \geq g_1$ .

To finalize the proof, note that the optimal value of the parameter  $g_0$  is  $\log(\gamma/(\gamma - \mathbb{E}^S \mathcal{R}(W)))$  and therefore since  $\gamma > 1$  and  $\mathbb{E}^S \mathcal{R}(W) \in [0, 1]$ , then  $g_0 \geq 0$ . Finally, letting  $c := e^{-g_0} \in (0, 1]$  and re-arranging the terms recovers the bound in the theorem.  $\blacksquare$

## A.2 Comparison between the fast-rate and mixed-rate bounds

Just by inspecting their equations, it is apparent that the proposed mixed-rate bound of Theorem 9 is tighter than those from Tolstikhin and Seldin (2013) and Rivasplata et al. (2019). However, it is not directly obvious that the presented fast-rate bound of Theorem 7 is tighter than Thiemann et al. (2017)'s Theorem 4. In fact, even Corollary 8 is tighter than this result.

To show this, we will show the stronger statement that  $f_{\text{fr}}(r, c) \leq f_{\text{th}}(r, c)$  for all  $r, c \geq 0$ , where

$$f_{\text{th}}(r, c) = \inf_{\lambda \in (0, 2)} \left\{ \frac{r}{1 - \frac{\lambda}{2}} + \frac{c}{\lambda(1 - \frac{\lambda}{2})} \right\} \text{ and } f_{\text{fr}}(r, c) = \inf_{\gamma > 2} \left\{ \gamma \log\left(\frac{\gamma}{\gamma - 1}\right) r + \gamma c \right\}.$$

If this holds, then Corollary 8 is tighter than Theorem 4 as enlarging the optimization set in  $f_{\text{fr}}(r, c)$  from  $\{\gamma > 2\}$  to  $\{\gamma > 1\}$  will only improve the bound.

Note that with the change of variable  $\gamma = (\lambda(1 - \lambda/2))^{-1}$ , if  $\lambda \in (0, 2)$ , then  $\gamma > 2$ . This way, we may re-write  $f_{\text{fr}}$  in terms of a minimization over  $\lambda \in (0, 2)$

$$f_{\text{fr}}(r, c) = \inf_{\lambda \in (0, 2)} \left\{ \frac{r}{\lambda(1 - \frac{\lambda}{2})} \log \frac{2}{(\lambda - 2)\lambda + 2} + \frac{c}{\lambda(1 - \frac{\lambda}{2})} \right\}.$$

Finally, noting that

$$\frac{1}{\lambda} \log \frac{2}{(\lambda - 2)\lambda + 2} \leq 1$$

for all  $\lambda \in (0, 2)$  completes the proof.

Similarly, it can also be shown that the mixed-rate bound from Theorem 9, which is itself a relaxation of the fast-rate bound of Corollary 8, is also tighter than the fast-rate bound from (Thiemann et al., 2017). In this case, we will show the stronger statement that  $f_{\text{mr}}(r, c) \leq f_{\text{th}}(r, c)$  for all  $r, c \geq 0$ , where

$$f_{\text{mr}}(r, c) = \inf_{\gamma > 2} \left\{ \frac{1}{2} \cdot \frac{2\gamma - 1}{\gamma - 1} r + \gamma c \right\}. \quad (13)$$

As above, as Theorem 9 is the closed-form expression obtained optimizing the equivalent of (13) on the larger set  $\{\gamma > 1\}$  (that is, optimizing (7)), showing this statement suffices.

Again, letting  $\gamma = (\lambda(1 - \lambda/2))^{-1}$ , if  $\lambda \in (0, 2)$  allows us to write  $f_{\text{mr}}$  in terms of a minimization over  $\lambda \in (0, 2)$

$$f_{\text{mr}}(r, c) = \inf_{\lambda \in (0, 2)} \left\{ \frac{1}{2} \cdot \frac{\lambda^2 - 2\lambda + 4}{\lambda^2 - 2\lambda + 2} r + \frac{c}{\lambda(1 - \frac{\lambda}{2})} \right\}.$$

Finally, noting that

$$\frac{1}{2} \cdot \frac{\lambda^2 - 2\lambda + 4}{\lambda^2 - 2\lambda + 2} \leq \frac{1}{1 - \frac{\lambda}{2}}$$

for all  $\lambda \in (0, 2)$  completes the proof.

### A.3 Example: PAC-Bayes with backprop using the fast and mixed-rate bounds

In Section 2, we mentioned that methods that use PAC-Bayes bounds to optimize the posterior, such as PAC-Bayes with backprop (Rivasplata et al., 2019; Pérez-Ortiz et al., 2021), could benefit from using the bounds from Theorems 7 and 9. In this subsection of the appendix, we provide an example showcasing that this is the case.

The PAC-Bayes with backprop method (Rivasplata et al., 2019; Pérez-Ortiz et al., 2021) considers a model  $\mathbf{n}_w$  parameterized by  $w \in \mathbb{R}^d$  and a prior distribution  $\mathbb{Q}_W$  over the parameters, for instance  $\mathbb{Q}_W = w_0 + \sigma_0 \mathcal{N}(0, I_d)$ . Then, the parameters are updated using stochastic gradient descent on the objective

$$\widehat{\mathcal{R}}(w, s) + f_{\text{bound}}(w; \mathbb{Q}_W),$$

where  $\widehat{\mathcal{R}}(w, s)$  is the empirical risk on the training data realization and  $f_{\text{bound}}(w; \mathbb{Q}_W)$  is extracted from a PAC-Bayes bound evaluated on the parameters  $w \in \mathbb{R}^d$  with prior  $\mathbb{Q}_W$ .

With an appropriate choice of the posterior  $\mathbb{P}_W^S$ , the bound function  $f_{\text{bound}}$  is calculable and the said posterior can be constructed, for instance  $\mathbb{P}_W^S(w) = w + \sigma_0 \mathcal{N}(0, I_d)$ . After the iterative procedure is completed, the empirical risk  $\mathbb{E}^{S=s} \widehat{\mathcal{R}}((, W), s)$  is bounded using the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) with a Monte Carlo estimate of the posterior parameters of  $m$  samples with confidence  $1 - \beta'$ , and the population risk is bounded also using the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) with the number of training samples  $n$  and a confidence  $1 - \beta$ , amounting for a total confidence of  $1 - (\beta' + \beta)$ . For more details, please check (Rivasplata et al., 2019; Pérez-Ortiz et al., 2021).

Using Thiemann et al. (2017)’s Theorem 4, Rivasplata et al. (2019)’s Theorem 5, or the classical McAllester (2003)’ bound (1) as an objective can be harmful since they penalize too harshly the dependence-confidence term dominated by the normalized dependency  $D(\mathbb{P}_W^S \| \mathbb{Q}_W) / n$ . Hence, SGD steers the parameters towards places too close to the prior, potentially avoiding other posteriors that achieve lower empirical error and have an overall better population risk. In this sense, it makes sense that bounds such as the proposed fast- and mixed-rate from Theorems 7 and 9 or the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002), with Reeb et al. (2018)’s gradients, would lead to said posteriors. This is verified in Table 1 for a convolutional network and the MNIST dataset. For the fast-rate bound from Theorem 7 and corollary 8, at each iteration the approximately optimal  $\gamma$  given after the theorem is employed, thus updating the posterior and the parameter alternately. We saw that the approximation of  $\gamma$  is good both by comparing the results of the final posterior in Table 1 and the coefficients of the empirical risk and the dependence-confidence term in Figure 1 with those obtained from the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) with Reeb et al. (2018, Appendix A)’s gradients. After a few iterations, once the empirical risk is small and Corollary 8 is a good approximation of Theorem 7, the gradients are close to each other.

**Remark 20** *Lotfi et al. (2022) obtain even tighter population risk certificates for networks on the MNIST dataset (11.6 %) considering a compression approach to the PAC Bayes bound from Catoni (2007). Therefore, their results could be tightened further using our strengthened version from Theorem 6. Nonetheless, the goal of this example is not to propose a method that obtains state-of-the-art certificates, but to showcase that the tightness of the tractable bounds in Section 2 can improve methods that employ PAC-Bayes bounds to find a suitable posterior.*

### A.3.1 EXPERIMENTAL DETAILS

All calculations were performed using the original code from PAC-Bayes with backprop: <https://github.com/mperezortiz/PBB>. The file modified to include our bounds and the hard-coded gradients from Reeb et al. (2018) is `bounds.py`. The convolutional network architecture consists of two convolutional layers with 32 and 64 filters respectively and a kernel size of 3. The last convolutional layer is followed by a max pooling layer with a kernel size of 2 and two linear layers with 128 and 10 nodes respectively. Between all layers there is a ReLU activation function.



Table 1: Population risk certificate, empirical risk, and normalized dependency of the posterior obtained with PAC-Bayes with Backprop (Rivasplata et al., 2019; Pérez-Ortiz et al., 2021) using Gaussian priors and different objectives. The best risk certificates are highlighted in bold face, and the second best is highlighted in italics. \*This refers to the normalized dependency  $D(\mathbb{P}_W^S \parallel \mathbb{Q}_W)/n$ . \*\*The gradients for the Seeger–Langford (Langford and Seeger, 2001; Seeger, 2002) bound are not calculated from the bound but hard-coded following Reeb et al. (2018, Appendix A).

Objective	Risk certificate	Empirical risk	Dependency*
Theorem 5 (Rivasplata et al., 2019)	0.20870	0.11372	0.03117
Theorem 4 (Thiemann et al., 2017)	0.21159	0.11053	0.03526
Equation (1) (McAllester, 2003)	0.23658	0.23658	0.02715
Corollary 8 [ours]	<b>0.17501</b>	0.07054	0.04649
Theorem 9 [ours]	<i>0.19763</i>	0.09214	0.04159
Theorem 2 (Seeger–Langford bound)**	<b>0.16922</b>	0.06701	0.04594

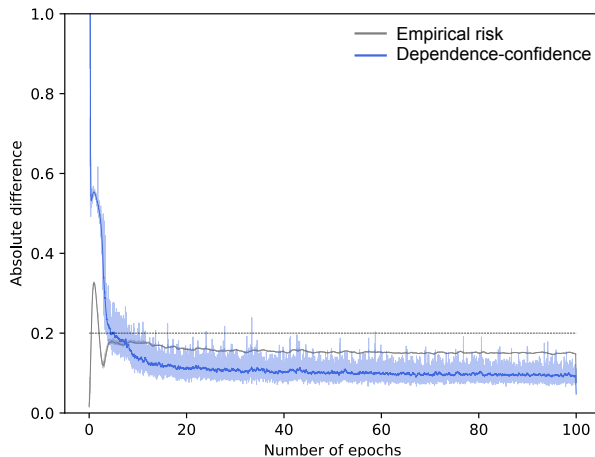


Figure 1: Absolute difference between the coefficients of the empirical risk (gray) and the dependence-confidence term (blue) of the gradients of the Seeger–Langford bound (Langford and Seeger, 2001; Seeger, 2002) from Reeb et al. (2018) and the fast-rate bound (Corollary 8) using the approximate optimal  $\gamma$ .

For all experiments, the standard deviation of the prior was  $\sigma_0 = 0.1$ . The learning rate was 0.01 for all experiments except for Rivasplata et al. (2019)’s Theorem 5 objective which was 0.005. The momentum was 0.99 for all objectives except for Thiemann et al. (2017)’s Theorem 4 which was 0.95. The number of Monte Carlo samples was  $m = 150,000$ , the minimum probability  $p_{\min}$  (see (Pérez-Ortiz et al., 2021) for the details) was  $10^{-5}$ , and the confidence parameters were  $\beta' = 0.01$  and  $\beta = 0.025$  respectively. The networks were trained for 100 epochs and a batch size of 250 to mimic the setting in (Pérez-Ortiz et al., 2021).

To find the hyper-parameters, we used the same grid search as Pérez-Ortiz et al. (2021). That is, the standard deviation of the prior was selected over  $\{0.005, 0.01, 0.02, 0.03, 0.04, 0.05, 0.1\}$ , the learning rate over  $\{0.001, 0.005, 0.01\}$ , and the momentum over  $\{0.95, 0.99\}$ . Therefore, the confidence parameters were updated to  $\beta' \leftarrow \beta'/42$  and  $\beta \leftarrow \beta/42$  respectively to comply with the union bound and maintain the guarantees.

All experiments were done on a TESLA V100 with 32GB of memory. Each full run takes approximately 110 hours with most of the time taken on the Monte Carlo sampling for the risk certificates calculation. For  $42 \cdot 5$  runs this amounts to approximately 23,100 hours which is around 32 months. Since the time was prohibitive for us, the hyper-parameter search was done without the Monte Carlo sampling, where each run took around 25 minutes amounting to a total of 87.5 hours or less than 4 days. Then, the final certificates were calculated using the full Monte Carlo sampling adding an extra 550 hours or around 23 days. In summary, the total amount of computing was approximately 27 days.

## Appendix B. Details of Section 3

This section of the appendix is devoted to providing alternative proofs and supplementary context and examples to the results from Section 3.

### B.1 Convex conjugates and their inverse

The convex conjugate is an important concept in convex analysis and in optimization theory. It is also particularly important to derive concentration inequalities through the Cramér–Chernoff method (Boucheron et al., 2013, Section 2.2) as shown in Section 3.

**Definition 21** *The convex conjugate (or just conjugate or Fenchel-Legendre’s dual) of a function  $\psi : \mathbb{R} \rightarrow \mathbb{R}$  is defined as*

$$\psi_*(x) := \sup_{\lambda \in \text{dom}(\psi)} \{\lambda x - \psi(\lambda)\}.$$

Specifically, when the function  $\psi$  is convex, the convex conjugate is also known as the *Legendre’s transform*, and when  $\psi$  represents or dominates a CGF as in Definition 10, it is known as the *Cramér’s transform*. We fall under both of these situations so the particular results for these transforms apply. A particularly important result is the following, which states an expression of the inverse of the convex conjugate of a smooth convex function. This result is used both to obtain the classical Chernoff’s inequality (Boucheron et al., 2013, Section 2.2) and its PAC-Bayes analogue from Theorem 12.

**Lemma 22 (Boucheron et al. (2013, Lemma 2.4))** *Let  $\psi$  be a convex and continuously differentiable function defined on  $[0, b)$  where  $0 < b \leq \infty$ . Assume that  $\psi(0) = \psi'(0) = 0$ . Then, the convex conjugate  $\psi_*$  is a non-negative convex and non-decreasing function on  $[0, \infty)$ . Moreover, for every  $y \geq 0$ , the set  $\{x \geq 0 : \psi_*(x) > y\}$  is non-empty and the generalized inverse of  $\psi_*$ , defined as  $\psi_*^{-1}(y) := \inf\{x \geq 0 : \psi_*(x) > y\}$  can also be written as*

$$\psi_*^{-1}(y) = \inf_{\lambda \in (0, b)} \left\{ \frac{y + \psi(\lambda)}{\lambda} \right\}.$$

## B.2 Proof of Theorem 11

Similarly to what we did in Appendix A, we first introduced Gibbs (1902, Theorem III on Chapter XI)' variational principle, which is a dual formulation to the Donsker and Varadhan (1975) Lemma 18.

**Lemma 23 (Gibbs' variational principle)** *Let  $\mathcal{X}$  be a measurable space,  $\mathbb{Q}$  be a probability measure on  $\mathcal{X}$ , and  $Y$  be a random variable distributed according to  $\mathbb{Q}$ . Further let  $g$  be a measurable function on  $\mathcal{X}$  such that  $\mathbb{E}[e^{g(Y)}] < \infty$  and  $\mathcal{P}_{\mathbb{Q}}(\mathcal{X})$  be the set of all probability measures  $\mathbb{P}$  on  $\mathcal{X}$  such that  $\mathbb{P} \ll \mathbb{Q}$ . Then,*

$$\log \mathbb{E} \left[ e^{g(Y)} \right] = \sup_{\mathbb{P} \in \mathcal{P}_{\mathbb{Q}}(\mathcal{X})} \left\{ \mathbb{E}[g(X)] - D(\mathbb{P} \parallel \mathbb{Q}) \right\}.$$

**Alternative proof of Theorem 11** Gibbs (1902)' variational principle from Lemma 23 states that for all measurable functions  $g$  such that  $e^g$  is  $\mathbb{Q}_W$  integrable

$$\log \mathbb{E}^S \left[ e^{g(W')} \right] = \sup_{\mathbb{P}_W^S \in \mathcal{P}_{\mathbb{Q}}(\mathcal{W})} \left\{ \mathbb{E}^S g(W) - D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) \right\} \text{ a.s.}, \quad (14)$$

where  $W'$  is distributed according to  $\mathbb{Q}_W$  and where  $\mathcal{P}_{\mathbb{Q}}(\mathcal{W})$  is the set of all measures on  $\mathcal{W}$  such that  $\mathbb{P}_W^S \ll \mathbb{Q}$  a.s.. In this case, let  $g(w; s) := \lambda(\mathcal{R}(w) - \widehat{\mathcal{R}}(w, s))$  for some  $\lambda \in (0, b/n)$ . The first term in the right hand side of (14) is directly  $\lambda \mathbb{E}^S[\mathcal{R}(W) - \widehat{\mathcal{R}}(W, S)]$ . For the term in the left hand side, one may employ Markov's inequality and Fubini's theorem to see that with probability no smaller than  $1 - \beta$

$$\mathbb{E}^S \left[ e^{g(W'; S)} \right] \leq \frac{1}{\beta} \cdot \mathbb{E} \left[ \mathbb{E}^{W'} \left[ e^{g(W'; S)} \right] \right].$$

Then, since  $\Lambda_{-\ell(w, Z)}(\lambda) \leq \psi(\lambda)$  for all  $w \in \mathcal{W}$  it holds that  $\Lambda_{\mathcal{R}(w) - \widehat{\mathcal{R}}(w, S)} \leq n\psi(\lambda/n)$ . Indeed,

$$\log \mathbb{E} \left[ e^{\lambda(\mathcal{R}(w) - \widehat{\mathcal{R}}(w, S))} \right] = \log \mathbb{E} \left[ e^{\frac{\lambda}{n} \sum_{i=1}^n (\mathbb{E}\ell(w, Z) - \ell(w, Z_i))} \right] \leq n\psi \left( \frac{\lambda}{n} \right).$$

Therefore, with probability no smaller than  $1 - \beta$

$$\log \mathbb{E}^S \left[ e^{\lambda(\mathcal{R}(W') - \widehat{\mathcal{R}}(W', S))} \right] \leq n\psi \left( \frac{\lambda}{n} \right) + \log \frac{1}{\beta}.$$

Combining the results for both terms together, for all  $\lambda \in (0, b/n)$ , with probability larger or equal than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \frac{1}{\lambda} \left[ D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) + \log \frac{1}{\beta} + n\psi \left( \frac{\lambda}{n} \right) \right]$$

holds *simultaneously* for all  $\mathbb{P}_W^S \in \mathcal{P}_{\mathbb{Q}_W}(\mathcal{W})$ , where the fact that it holds uniformly for all the Markov kernels  $\mathbb{P}_W^S$  comes from the supremum. Since the result holds for all  $\lambda \in (0, b/n)$ , performing the substitution  $\lambda \leftarrow n\lambda$  completes the proof.  $\blacksquare$

### B.3 The optimization of $\lambda$ in (10)

Recall that in the proof of Theorem 12 we considered the event  $\mathcal{B}_\lambda$  to be the complement of the event in (8) such that  $\mathbb{P}[\mathcal{B}_\lambda] \leq \beta$  for all  $\lambda \in (0, b)$ . This event is parameterized with  $\lambda$  and, given the event  $\mathcal{E}_k = \{\lceil D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) \rceil = k\}$ , with probability no more than  $\mathbb{P}[\mathcal{B}_\lambda | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \frac{1}{\lambda} \left[ \frac{k + \log \frac{1}{\beta}}{n} + \psi(\lambda) \right] \quad (10)$$

for all  $\lambda \in (0, b)$ . Then, after optimizing the parameter  $\lambda$  on the right hand side of (10) using Lemma 22, we considered the event  $\mathcal{B}_{\lambda_k}$  resulting of that optimization such that  $\mathbb{P}[\mathcal{B}_{\lambda_k}] \leq \beta$ . This notation is imprecise, the infimum in (10) is attained either by a  $\lambda_k \in (0, b)$  or by letting  $\lambda \rightarrow b$ . It will never be attained when  $\lambda \rightarrow 0$  as  $\psi(0) = 0$  and the term inside the infimum goes to  $\infty$  when  $\lambda \rightarrow 0$ . In the case where the infimum is attained by letting  $\lambda \rightarrow b$ , by continuity, the desired inequality (11) still holds and the event described by  $\lim_{\lambda \rightarrow b} \mathcal{B}_\lambda$  is still such that  $\mathbb{P}[\lim_{\lambda \rightarrow b} \mathcal{B}_\lambda] \leq \beta$ . We hide these details from the main text for clarity of exposition.

### B.4 PAC-Bayes bounds for different loss tail behaviors

Theorem 12 describes different high-probability PAC-Bayes bounds for different tail behaviors of the loss  $\ell(w, Z)$ . The next corollary collects some of the most common tail behaviors and their resulting PAC-Bayes bound.

**Corollary 24** *Consider a training set  $S$  with  $n$  samples. Let  $\mathbb{Q}_W$  be any prior independent of  $S$  and define  $\mathfrak{C}_{n,\beta,S} := \frac{1}{n} (D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) + \log \frac{1}{\beta})$  and the event  $\mathcal{E} := \{D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) \leq n\}$ . Then:*

1. *if the loss  $\ell$  has a bounded range  $[a, b]$ , where  $-\infty < a \leq b < \infty$ , then with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \sqrt{(b-a)^2 \mathfrak{C}_{2n,\beta/en,S}}$$

*holds simultaneously for every posterior  $\mathbb{P}_W^S$ ;*

2. *if the loss  $\ell(w, Z)$  is  $\sigma^2$ -subgaussian for all hypotheses  $w \in \mathcal{W}$ , then with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbf{1}_{\mathcal{E}} \cdot \left[ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \sqrt{2\sigma^2 \mathfrak{C}_{n,\beta/en,S}} \right] + \mathbf{1}_{\mathcal{E}^c} \cdot \text{ess sup } \mathbb{E}^S \mathcal{R}(W)$$

*holds simultaneously for every posterior  $\mathbb{P}_W^S$ ;*

3. *if the loss  $\ell(w, Z)$  is  $(\sigma^2, c)$ -subgamma for all hypotheses  $w \in \mathcal{W}$ , then with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbf{1}_{\mathcal{E}} \cdot \left[ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \sqrt{2\sigma^2 \mathfrak{C}_{n,\beta/en,S} + c \mathfrak{C}_{n,\beta/en,S}} \right] + \mathbf{1}_{\mathcal{E}^c} \cdot \text{ess sup } \mathbb{E}^S \mathcal{R}(W)$$

*holds simultaneously for every posterior  $\mathbb{P}_W^S$ ;*

4. and if  $\ell(w, Z)$  is  $(\sigma^2, c)$ -subexponential<sup>9</sup> for all hypotheses  $w \in \mathcal{W}$ , then with probability no smaller than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbf{1}_{\mathcal{E}} \cdot \left[ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \mathbf{1}_{\mathcal{F}} \cdot \sqrt{2\sigma^2 \mathfrak{C}_{n, \beta/en, S}} + \mathbf{1}_{\mathcal{F}^c} \cdot (c+1) \mathfrak{C}_{n, \beta/en, S} \right] + \mathbf{1}_{\mathcal{E}^c} \cdot \text{ess sup } \mathbb{E}^S \mathcal{R}(W)$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ , where  $\mathcal{F}$  is the event  $\mathcal{F} := \{D(\mathbb{P}_W^S \| \mathbb{Q}_W) \leq \frac{n\sigma^2}{2c} - \log \frac{\epsilon}{\beta}\}$ .

**Proof** We may prove each point individually:

- Point 2 follows by noting that for  $\sigma^2$ -subgaussian random variables  $\psi(\lambda) = \lambda^2 \sigma^2 / 2$  and therefore  $\psi_*^{-1}(y) = \sqrt{2\sigma^2 y}$ .
- Point 1 follows by noting that if a random variable is bounded in  $[a, b]$ , then it is  $(b-a)^2/2$ -subgaussian. Then, as hinted in Section 3.2.2, we may consider the more lenient event  $\mathcal{E}' := \{D(\mathbb{P}_W^S \| \mathbb{Q}_W) \leq 2n\}$  and use that  $\text{ess sup } \mathbb{E}^S \mathcal{R}(W) \leq b - a \leq \sqrt{(b-a)^2 \mathfrak{C}_{2n, \beta/en, S}}$ .
- Point 3 follows by noting that for  $(\sigma^2, c)$ -subgamma random variables  $\psi(\lambda) = \lambda^2 \sigma^2 / 2(1-c\lambda)$  for all  $\lambda \in (0, 1/c)$  and therefore  $\psi_*^{-1}(y) = \sqrt{2\sigma^2 y} + cy$  (Boucheron et al., 2013, Section 2.4).
- Finally, Point 4 follows by noting that for  $(\sigma^2, c)$ -subexponential ranom objects  $\psi(\lambda) = \lambda^2 \sigma^2 / 2$  for all  $\lambda \in (0, 1/c)$  and therefore

$$\psi_*^{-1}(y) = \begin{cases} \sqrt{2\sigma^2 y} & \text{if } \lambda = \sqrt{2y/\sigma^2} \leq 1/c \\ cy + \sigma^2/2c^2 & \text{otherwise} \end{cases}.$$

The condition for the first case may be rewritten as  $y \leq \sigma^2/2c^2$  and similarly the condition for the second case as  $y > \sigma^2/2c^2$ . Hence, we have the inequality

$$\psi_*^{-1}(y) \leq \begin{cases} \sqrt{2\sigma^2 y} & \text{if } y \leq \sigma^2/2c^2 \\ (c+1)y & \text{otherwise} \end{cases}.$$

■

### B.5 A parameter-free version of Wang et al. (2015) and Haddouche and Guedj (2023a, Theorem 2.1)'s PAC-Bayes bound on martingales

Wang et al. (2015) and Haddouche and Guedj (2023a) investigate the setting where the dataset  $S$  is considered to be a sequence  $S^* := (Z_i)_{i \geq 1}$  such that  $Z_i \in \mathcal{Z}$ , but where there is

9. Here, we are considering the subexponential characterization of random variables from Wainwright (2019, Theorem 2.13), and not the one given by Boucheron et al. (2013, Exercice 2.22).

no restriction in the distribution of the samples  $Z_i$ , that is, every sample  $Z_i$  can depend on all the previous ones. For every  $n$ , they let  $S_n := (Z_1, \dots, Z_n)$  be the restriction of  $S^*$  to its first  $n$  points. Then, they consider the sequence of  $\sigma$ -algebras  $(\mathcal{F}_i)_{i \geq 1}$  to be a filtration adapted to  $S^*$ , for instance  $\mathcal{F}_i = \sigma(Z_1, \dots, Z_i)$ . Finally, they consider a martingale difference sequence  $(X_i(S_i, w))_{i \geq 1}$  indexed by a hypothesis  $w \in \mathcal{W}$  so that  $\mathbb{E}^{\mathcal{F}_{i-1}} X_i(S_i, w) = 0$  for all  $w \in \mathcal{W}$ . For instance, let  $Y_0 = \sum_{i=1}^n \mathbb{E} \ell(w, Z_i)$  and  $Y_i(S_i, w) = \sum_{i=1}^n \mathbb{E}^{\mathcal{F}_i} \ell(w, Z_i)$  for all  $i \geq 1$ , then  $X_i(S_i, w) = Y_i - Y_{i-1}$ . Finally, for all  $w \in \mathcal{W}$ , they define the martingale  $M_n(w) := \sum_{i=1}^n X_i(S_i, w)$  and follow Bercu and Touati (2008) to also define

$$[M]_n(w) := \sum_{i=1}^n X_i(S_i, w)^2 \text{ and } \langle M \rangle_n(w) := \mathbb{E}^{\mathcal{F}_{i-1}} X_i(S_i, w)^2,$$

where  $[M]_n(w)$  acts as an empirical variance term and  $\langle M \rangle_n(w)$  as its theoretical counterpart (Haddouche and Guedj, 2023a).

Then, their main anytime-valid bound for martingales is the following.

**Theorem 25 (Wang et al. (2015, Theorem 2.4))** *Let  $\mathbb{Q}_W$  be any prior independent of  $S_n$  and  $(M_n(w))_{n \geq 1}$  be any collection of martingales indexed by  $w \in \mathcal{W}$ . Then, for all  $\lambda > 0$ , all  $\beta \in (0, 1)$ , and simultaneously for all  $n \geq 1$ , with probability no smaller than  $1 - \beta$*

$$|\mathbb{E}^{S_n} M_n(W)| \leq \frac{D(\mathbb{P}_W^{S_n} \parallel \mathbb{Q}_W) + \log \frac{2(n+1)^2}{\beta}}{\lambda} + \frac{\lambda}{6} \cdot \mathbb{E}^{S_n} [[M]_n(W) + 2\langle M \rangle_n(W)] \quad (15)$$

holds simultaneously for every posterior  $\mathbb{P}_W^{S_n}$ .

**Theorem 26 (Haddouche and Guedj (2023a, Theorem 2.1))** *Let  $\mathbb{Q}_W$  be any prior independent of  $S_n$  and  $(M_n(w))_{n \geq 1}$  be any collection of martingales indexed by  $w \in \mathcal{W}$ . Then, for all  $\lambda > 0$ , all  $\beta \in (0, 1)$ , and simultaneously for all  $n \geq 1$ , with probability no smaller than  $1 - \beta$*

$$|\mathbb{E}^{S_n} M_n(W)| \leq \frac{D(\mathbb{P}_W^{S_n} \parallel \mathbb{Q}_W) + \log \frac{2}{\beta}}{\lambda} + \frac{\lambda}{2} \cdot \mathbb{E}^{S_n} [[M]_n(W) + \langle M \rangle_n(W)]$$

holds simultaneously for every posterior  $\mathbb{P}_W^{S_n}$ .

In what follows, we will focus on the result from Wang et al. (2015) as it has the smaller constants. Taking a closer look at Theorem 25, we realize it has a similar shape to Theorem 11 for the particular case when the loss is subgaussian, where the role of the subgaussian parameter is taken by the sum of the “variance” terms  $[M]_n(W) + 2\langle M \rangle_n(W)$ . Therefore, it appears we may directly employ the technique to derive the Chernoff analogue from the proof of Theorem 12. However, one needs to take into account the fact that the “optimal” parameter  $\lambda$  now depends on this “variance” terms, which are also dependent on the training set  $S_n$  and on the number of samples  $n$ .

To optimize the bound from Theorem 25 we will then proceed in two steps. The first step is to optimize the parameter  $\lambda$  for a *fixed* number of samples  $n$  in a similar fashion to Theorem 12, which results in Theorem 27. Then, the second step is to extend this result to

an anytime valid bound using Theorem 16 at a cost in the dependence-confidence term of  $\mathcal{O}(\log n/n)$ .

For the first step, define the event  $\mathcal{B}_{n,\lambda}$  as the complement of the event in (15) for a *fixed* number of samples  $n$ . Then, we can proceed similarly to the proof of Theorem 12 noticing that, for each number of samples  $n$ , the complement of the event

$$\mathcal{E}_n := \left\{ \mathbb{E}^{S_n} [[M]_n(W) + 2\langle M \rangle_n(W)] D(\mathbb{P}_W^S \|\mathbb{W}_W) \leq n^2 \right\}$$

is uninteresting as the bound is non-vanishing given  $\mathcal{E}_n^c$ . This produces the following PAC-Bayes bound for a fixed number of samples  $n$ .

**Theorem 27 (Parameter-free bound on martingales)** *Let  $\mathbb{Q}_W$  be any prior independent of  $S_n$  and  $(M_n(w))_{n \geq 1}$  be any collection of martingales indexed by  $w \in \mathcal{W}$ . Further, define  $\xi'(n) := 2en(n+1)^2 \log(en) \leq 2e(n+1)^3$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$*

$$\begin{aligned} |\mathbb{E}^{S_n} M_n(W)| &\leq \mathbf{1}_{\mathcal{E}_n} \cdot \frac{2}{\sqrt{6}} \cdot \sqrt{\mathbb{E}^{S_n} [[M]_n(W) + 2\langle M \rangle_n(W) + 1] \left( D(\mathbb{P}_W^S \|\mathbb{Q}_W) + \log \frac{\xi'(n)}{\beta_n} \right)} \\ &\quad + \mathbf{1}_{\mathcal{E}_n^c} \text{ess sup } |\mathbb{E}^{S_n} M_n(W)| \end{aligned}$$

holds simultaneously for all posteriors  $\mathbb{P}_W^{S_n}$ , where  $\mathcal{E}_n$  is the event  $\mathcal{E}_n := \left\{ \mathbb{E}^{S_n} [[M]_n(W) + 2\langle M \rangle_n(W)] D(\mathbb{P}_W^S \|\mathbb{W}_W) \leq n^2 \right\}$ .

**Proof** Consider a fixed number of samples  $n$ . Let  $\mathcal{B}_{n,\lambda}$  be the complement of the event in (15) such that  $\mathbb{P}[\mathcal{B}_{n,\lambda}] < \beta$  and consider the sub-events

$$\begin{aligned} \mathcal{E}_{n,1,l} &:= \left\{ D(\mathbb{P}_W^{S_n} \|\mathbb{Q}_W) \leq 1 \text{ and } \lceil \mathbb{E}^{S_n} [[M]_n(W) + 2\langle M \rangle_n(W)] \rceil = l \right\}, \\ \mathcal{E}_{n,k,1} &:= \left\{ \lceil D(\mathbb{P}_W^{S_n} \|\mathbb{Q}_W) \rceil = k \text{ and } \mathbb{E}^{S_n} [[M]_n(W) + 2\langle M \rangle_n(W)] \leq 1 \right\}, \text{ and} \\ \mathcal{E}_{n,k,l} &:= \left\{ \lceil D(\mathbb{P}_W^{S_n} \|\mathbb{Q}_W) \rceil = k \text{ and } \lceil \mathbb{E}^{S_n} [[M]_n(W) + 2\langle M \rangle_n(W)] \rceil = l \right\}, \end{aligned}$$

for all  $k, l = 2, \dots, n^2$  such that  $kl \leq n^2$ , which form a covering of  $\mathcal{E}_n$ . Furthermore, define  $\mathcal{K} := \{(k, l) : 1 \leq kl \leq n^2 \text{ and } \mathbb{P}[\mathcal{E}_{n,k,l}] > 0\}$ . For all  $(k, l) \in \mathcal{K}$ , given the event  $\mathcal{E}_{n,k,l}$ , with probability no more than  $\mathbb{P}[\mathcal{B}_{n,\lambda} | \mathcal{E}_{n,k,l}]$ , there exists some posterior  $\mathbb{P}_W^{S_n}$  such that

$$|\mathbb{E}^{S_n} M_n(W)| > \frac{k + \log \frac{2(n+1)^2}{\beta}}{\lambda} + \frac{\lambda}{6} \cdot l. \quad (16)$$

for all  $\lambda \in (0, b)$ . The parameter that optimizes the right hand side of (16) is

$$\lambda = \lambda_{k,l} = \sqrt{\frac{6}{l} \left( k + \log \frac{2(n+1)^2}{\beta} \right)}.$$

Substituting the optimal  $\lambda_{k,l}$  and using that  $k \leq D(\mathbb{P}_W^{S_n} \|\mathbb{Q}_W) + 1$  and  $l \leq \mathbb{E}^{S_n} [[M]_n(W) + 2\langle M \rangle_n(W) + 1]$  yields that, given the event  $\mathcal{E}_{n,k,l}$ , with probability smaller or equal than  $\mathbb{P}[\mathcal{B}_{n,\lambda_{k,l}} | \mathcal{E}_{n,k,l}]$ , there exists some posterior  $\mathbb{P}_W^{S_n}$  such that

$$|\mathbb{E}^{S_n} M_n(W)| > \frac{2}{\sqrt{6}} \cdot \sqrt{\mathbb{E}^{S_n} [[M]_n(W) + \langle M \rangle_n(W) + 1] \left( D(\mathbb{P}_W^{S_n} \|\mathbb{Q}_W) + \log \frac{2e(n+1)^2}{\beta} \right)}.$$

Now, define  $\mathcal{B}'_n$  as the event stating that there exists some posterior  $\mathbb{P}_W^{S_n}$  such that

$$|\mathbb{E}^{S_n} M_n(W)| > \mathbf{1}_{\mathcal{E}} \cdot \frac{2}{\sqrt{6}} \cdot \sqrt{\mathbb{E}^{S_n} [[M]_n(W) + 2\langle M \rangle_n(W) + 1] \left( D(\mathbb{P}_W^{S_n} \|\mathbb{Q}_W) + \log \frac{2e(n+1)^2}{\beta} \right)} + \mathbf{1}_{\mathcal{E}} \cdot \text{ess sup } |\mathbb{E}^{S_n} M_n(W)|,$$

where  $\mathbb{P}[\mathcal{B}'_n | \mathcal{E}_{n,k,l}] \mathbb{P}[\mathcal{E}_{n,k,l}] \leq \mathbb{P}[\mathcal{B}_{n,\lambda_{k,l}} | \mathcal{E}_{n,k,l}] \mathbb{P}[\mathcal{E}_{n,k,l}] \leq \mathbb{P}[\mathcal{B}_{n,\lambda_{k,l}}] < \beta$  for all  $(k,l) \in \mathcal{K}$ , and where  $\mathbb{P}[\mathcal{B}'_n \cap \mathcal{E}_n^c] = 0$  by the definition of the essential supremum. Therefore, the probability of  $\mathcal{B}'_n$  is bounded as

$$\mathbb{P}[\mathcal{B}'_n] = \sum_{(k,l) \in \mathcal{K}} \mathbb{P}[\mathcal{B}'_n | \mathcal{E}_{n,k,l}] \mathbb{P}[\mathcal{E}_{n,k,l}] + \mathbb{P}[\mathcal{B}'_n \cap \mathcal{E}_n^c] < n(1 + \log n)\beta = n \log(en)\beta.$$

Finally, let  $\beta_n = n \log(en)\beta$  so that, with probability no larger than  $\beta_n$ , there exists some posterior  $\mathbb{P}_W^{S_n}$  such that

$$|\mathbb{E}^{S_n} M_n(W)| > \mathbf{1}_{\mathcal{E}} \cdot \frac{2}{\sqrt{6}} \cdot \sqrt{\mathbb{E}^{S_n} [[M]_n(W) + 2\langle M \rangle_n(W) + 1] \left( D(\mathbb{P}_W^{S_n} \|\mathbb{Q}_W) + \log \frac{2en(n+1)^2 \log(en)}{\beta_n} \right)} + \mathbf{1}_{\mathcal{E}} \cdot \text{ess sup } |\mathbb{E}^{S_n} M_n(W)|.$$

Finally, the substitution  $\beta_n \leftarrow \beta$  completes the proof.  $\blacksquare$

This technique can be extended to the corollary bound of Haddouche and Guedj (2023a) for batch learning with i.i.d. data yielding Theorem 13, where we write  $S_n = S$  to simplify the reading in the main text. Note again that we are using the particularization of Haddouche and Guedj (2023a) with the constants from Wang et al. (2015).

Finally, for the second step, Theorems 13 and 27 can be converted back to anytime-valid bounds using Theorem 16. The resulting bound is exactly the same substituting  $\log \xi'(n)$  for  $\log \xi''(n)$ , where  $\xi''(n) := e\pi^2(n+1)^2 n^3 \log(en)/3$ .

In case that one desires to have a bound without a  $\log n$  term, one may consider employing the technique outlined for Theorem 14 with Haddouche and Guedj (2023a)'s Theorem 26 instead of Wang et al. (2015)'s Theorem 25.

## B.6 PAC-Bayes bounds with a smaller union bound cost

As discussed in Section 3, the union bound cost of the PAC-Bayes bounds developed above can be improved at the cost of a multiplicative factor of  $e$  to the relative entropy. Below, we present the parallel of Theorem 12 with this improved union bound cost, but extending Corollary 24 and Theorem 13, 27, and 31 follows analogously almost verbatim.



**Theorem 28** Consider a loss function  $\ell$  with a bounded CGF in the sense of Definition 10. Let  $\mathbb{Q}_W$  be any prior independent of  $S$  and define the event  $\mathcal{E} = \{\mathbb{D}(\mathbb{P}_W^S \parallel \mathbb{Q}_W) \leq n\}$ . Then, for every  $\beta \in (0, 1)$  with probability no smaller than  $1 - \beta$

$$\begin{aligned} \mathbb{E}^S \mathcal{R}(W) &\leq \mathbb{1}_{\mathcal{E}} \cdot \left[ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{e \max\{\mathbb{D}(\mathbb{P}_W^S \parallel \mathbb{Q}_W), 1\} + \log \frac{2 + \log n}{\beta}}{n} \right) \right] \\ &\quad + \mathbb{1}_{\mathcal{E}^c} \cdot \text{ess sup } \mathbb{E}^S \mathcal{R}(W) \end{aligned}$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

**Proof** Let  $\mathcal{B}_\lambda$  be the complement of the event in (8) such that  $\mathbb{P}[\mathcal{B}_\lambda] < \beta$  and consider the sub-events  $\mathcal{E}_0 := \{\mathbb{D}(\mathbb{P}_W^S \parallel \mathbb{Q}_W) \in [0, 1]\}$ ,  $\mathcal{E}_1 := \{\mathbb{D}(\mathbb{P}_W^S \parallel \mathbb{Q}_W) \in (1, e]\}$ , and  $\mathcal{E}_k := \{\mathbb{D}(\mathbb{P}_W^S \parallel \mathbb{Q}_W) \in (e^{k-1}, e^k]\}$  for all  $k = 2, \dots, n$ , which form a covering of the event  $\mathcal{E} := \{\mathbb{D}(\mathbb{P}_W^S \parallel \mathbb{Q}_W) \leq n\}$ . Furthermore, define  $\mathcal{K} := \{k \in \mathbb{N} \cup \{0\} : 0 \leq k \leq n \text{ and } \mathbb{P}[\mathcal{E}_k] > 0\}$ . For all  $k \in \mathcal{K} \setminus \{0\}$ , given the event  $\mathcal{E}_k$ , with probability no more than  $\mathbb{P}[\mathcal{B}_\lambda | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \frac{1}{\lambda} \left[ \frac{e^k + \log \frac{1}{\beta}}{n} + \psi(\lambda) \right], \quad (17)$$

for all  $\lambda \in (0, b)$ . The right hand side of (17) can be minimized with respect to  $\lambda$  *independently of the training set*  $S$ . Let  $\mathcal{B}_{\lambda_k}$  be the event resulting from this minimization and note that  $\mathbb{P}[\mathcal{B}_{\lambda_k}] < \beta$ . According to (Boucheron et al., 2013, Lemma 2.4), this ensures that, with probability no more than  $\mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{e^k + \log \frac{1}{\beta}}{n} \right),$$

where  $\psi_*$  is the convex conjugate of  $\psi$  and where  $\psi_*^{-1}$  is a non-decreasing concave function. Given  $\mathcal{E}_k$ , since  $e^k < e\mathbb{D}(\mathbb{P}_W^S \parallel \mathbb{Q}_W)$ , with probability no larger than  $\mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{e\mathbb{D}(\mathbb{P}_W^S \parallel \mathbb{Q}_W) + \log \frac{1}{\beta}}{n} \right).$$

Now, define  $\mathcal{B}'$  as the event stating that there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{1}_{\mathcal{E}} \cdot \left[ \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{e \max\{\mathbb{D}(\mathbb{P}_W^S \parallel \mathbb{Q}_W), 1\} + \log \frac{e}{\beta}}{n} \right) \right] + \mathbb{1}_{\mathcal{E}^c} \cdot \text{ess sup } \mathbb{E}^S \mathcal{R}(W)$$

where  $\mathbb{P}[\mathcal{B}' | \mathcal{E}_k] \mathbb{P}[\mathcal{E}_k] \leq \mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k] \mathbb{P}[\mathcal{E}_k] \leq \mathbb{P}[\mathcal{B}_{\lambda_k}] \leq \beta$  for all  $k \in \mathcal{K}$ , and where  $\mathbb{P}[\mathcal{B}' \cap \mathcal{E}^c] = 0$  by the definition of the essential supremum. Note that, if  $\{0\} \in \mathcal{K}$ , the case for  $k = 0$  is handled by the addition of the maximum  $\max\{\mathbb{D}(\mathbb{P}_W^S \parallel \mathbb{Q}_W), 1\}$  to the equation defining the event  $\mathcal{B}'$ . Therefore, the probability of  $\mathcal{B}'$  is bounded as

$$\mathbb{P}[\mathcal{B}'] = \sum_{k \in \mathcal{K}} \mathbb{P}[\mathcal{B}' | \mathcal{E}_k] \mathbb{P}[\mathcal{E}_k] + \mathbb{P}[\mathcal{B}' \cap \mathcal{E}^c] < (2 + \log n)\beta.$$

Finally, the substitution  $\beta \leftarrow \beta / (2 + \log n)$  completes the proof.  $\blacksquare$

### B.7 PAC-Bayes bounds without an uninteresting event

As discussed in Section 3, the presented approach to find parameter-free PAC-Bayes bounds can be extended to the case where no event is considered uninteresting. Below, we present the parallel of Theorem 12 with this consideration. However, extending Corollary 24 and Theorems 13, 27 and 31 and 28 follows analogously almost verbatim. The main important considerations are that to extend Theorems 13 and 27 one needs to consider a double sum  $\sum_{k=1}^{\infty} \sum_{l=1}^{\infty} \beta_{k,l}$  and that to extend the results with the log log cost from Appendix B.6 one needs to separate the space with a geometric grid.

**Proof of Theorem 14** Consider the sub-events  $\mathcal{E}_1 := \{D(\mathbb{P}_W^S \parallel \mathbb{Q}) \leq 1\}$  and  $\mathcal{E}_k := \{[D(\mathbb{P}_W^S \parallel \mathbb{Q})] = k\}$  for all  $k \geq 2 \in \mathbb{N}$ , which form a covering of the events' space. Furthermore, define  $\mathcal{K} := \{k \in \mathbb{N} : 1 \leq k \text{ and } \mathbb{P}[\mathcal{E}_k] > 0\}$ . For all  $k \in \mathcal{K}$ , consider the event  $\mathcal{B}_{\lambda_k}$  to be the complement of the event in (8) for a given parameter  $\beta_k$  such that  $\mathbb{P}[\mathcal{B}_{\lambda_k}] < \beta_k$ . Then, given the event  $\mathcal{E}_k$ , with probability no more than  $\mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \frac{1}{\lambda_k} \left[ \frac{k + \log \frac{1}{\beta_k}}{n} + \psi(\lambda_k) \right], \quad (18)$$

for all  $\lambda_k \in (0, b)$ . The right hand side of (18) can be minimized with respect to  $\lambda$  *independently of the training set*  $S$ . Let  $\mathcal{B}_{\lambda_k}$  be the event resulting from this minimization and note that  $\mathbb{P}[\mathcal{B}_{\lambda_k}] \leq \beta_k$ . According to (Boucheron et al., 2013, Lemma 2.4), this ensures that with probability no more than  $\mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{k + \log \frac{1}{\beta_k}}{n} \right),$$

where  $\psi_*$  is the convex conjugate of  $\psi$  and where  $\psi_*^{-1}$  is a non-decreasing concave function. Now, let  $\beta_k = \frac{\beta}{k^2}$ . Given  $\mathcal{E}_k$ , since  $k < D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) + 1$ , with probability no larger than  $\mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^S \mathcal{R}(W) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \psi_*^{-1} \left( \frac{D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) + 1 + \log \frac{(D(\mathbb{P}_W^S \parallel \mathbb{Q}_W) + 1)^2}{\beta}}{n} \right). \quad (19)$$

Now, define  $\mathcal{B}'$  as the event described in (19), where  $\mathbb{P}[\mathcal{B}' | \mathcal{E}_k] \mathbb{P}[\mathcal{E}_k] \leq \mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k] \mathbb{P}[\mathcal{E}_k] \leq \mathbb{P}[\mathcal{B}_{\lambda_k}] \leq \beta_k = \frac{\beta}{k^2}$  for all  $k \in \mathcal{K}$ . Therefore, the probability of  $\mathcal{B}'$  is bounded as

$$\mathbb{P}[\mathcal{B}'] = \sum_{k \in \mathcal{K}} \mathbb{P}[\mathcal{B}' | \mathcal{E}_k] \mathbb{P}[\mathcal{E}_k] < \sum_{k=1}^{\infty} \frac{\beta}{k^2} = \frac{\pi^2}{6} \cdot \beta.$$

Finally, the substitution  $\beta \leftarrow 6\beta/\pi^2$  completes the proof.  $\blacksquare$

### B.8 Example: Recovering a PAC-Bayes bound on the randomized subsample setting

In this section of the Appendix, as a further application of the technique devised to prove Theorem 12 to obtain bounds “in probability”, we recover Hellström and Durisi (2020)’s PAC-Bayes

bound in the *randomized subsample setting*. This setting was introduced by Steinke and Zakyntinou (2020) motivated by the fact that the dependency measure  $D(\mathbb{P}_{W|S} \|\mathbb{Q}_W)$  can be large, or even infinite, in situations where algorithms generalize (Bassily et al., 2018; Livni and Moran, 2020).

In the randomized subsample setting, it is considered that the training dataset  $S$  is obtained through the following mechanism. First, a super sample

$$\tilde{S} := \begin{pmatrix} \tilde{Z}_{1,1} & \tilde{Z}_{2,1} & \cdots & \tilde{Z}_{n,1} \\ \tilde{Z}_{1,2} & \tilde{Z}_{2,2} & \cdots & \tilde{Z}_{n,2} \end{pmatrix}^\top$$

of  $2n$  i.i.d. instances is obtained by sampling from the distribution  $\mathbb{P}_Z$ . Unused samples, or *ghost samples*  $S_{\text{ghost}} = \tilde{S} \setminus S$ , are virtual and only exist for the purpose of the analysis. Then, an independent sequence of indices  $U := (U_1, \dots, U_n)$  distributed as  $\mathbb{P}_{U_i}[1] = \mathbb{P}_{U_i}[2] = 1/2$  is generated. Finally, the training set is sub-sampled from the superset so that  $Z_i = \tilde{Z}_{i,U_i}$ . Under this setting, Steinke and Zakyntinou (2020) and Grunwald et al. (2021) derive PAC-Bayes and MAC-Bayes bounds, where the dependence measure is now the relative entropy  $D(\mathbb{P}_W^S \|\mathbb{Q}_W^{\tilde{S}})$  of the posterior  $\mathbb{P}_W^S$  with respect to a *data-dependent* prior  $\mathbb{Q}_W^{\tilde{S}}$  that has access to the supersample  $\tilde{S}$  but not to the indices  $U$ .

Bounds depending on this “conditional” measure are promising. Steinke and Zakyntinou (2020) and Grunwald et al. (2021) showed that for finite VC dimension and compression schemes both MAC- and PAC-Bayes bounds can be obtained; Haghifam et al. (2021) showed that a highly related measure provides a sharp characterization of the population risk in the realizable setting for 0–1 losses; and Hellström and Durisi (2022) made a similar remark for both PAC- and MAC-Bayes bounds for classes with finite Natarajan dimension (Natarajan, 1989). However, recently Haghifam et al. (2023) showed that there are still simple algorithms that generalize but where this conditional measure is high.

In this setting, the canonical PAC-Bayes bound is given by Hellström and Durisi (2021).

**Theorem 29 (Hellström and Durisi (2020, 2021, Theorem 3 and Corollary 6))**

Consider a bounded loss function  $\ell : \mathcal{W} \times \mathcal{Z} \rightarrow [a, b]$ . Let  $\mathbb{Q}_W^{\tilde{S}}$  be any Markov kernel on  $\mathcal{W}$  with access to the supersample  $\tilde{S}$  but not to the indices  $U$ . Then, for all  $\beta \in (0, 1)$ , with probability no less than  $1 - \beta$

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{E}^S \hat{\mathcal{R}}(W, S) + \sqrt{\frac{2(b-a)^2 (D(\mathbb{P}_W^S \|\mathbb{Q}_W^{\tilde{S}}) + \log \frac{2\sqrt{n}}{\beta})}{n-1}} + \sqrt{\frac{(b-a)^2 \log \frac{4}{\beta}}{2n}}$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ , where the conditioning is written on  $S$  and not on  $\tilde{S}$  and  $U$  since  $\mathbb{P}_W^{\tilde{S}, U} = \mathbb{P}_W^S$  a.s..

To obtain the PAC-Bayes bound from Theorem 29, Hellström and Durisi (2021) use a specific property of subgaussian random variables (Wainwright, 2019, Theorem 2.6). In the randomized subsample setting, where the loss is assumed to be bounded, this results in a general bound as bounded random variables are subgaussian. Still, to illustrate the technique from Theorem 12, we show how an equivalent result is obtained based on Hellström and Durisi (2020)’s original procedure.

Using Hellström and Durisi (2020, Theorem 4)’s exponential inequality, which independently re-discovered (Zhang, 2006, Lemma 2.1), and following the steps of (Hellström and Durisi, 2020, Corolary 2) yields the next result.

**Lemma 30** *Consider a bounded loss function  $\ell : \mathcal{W} \times \mathcal{Z} \rightarrow [a, b]$ . Let  $\mathbb{Q}_W^{\tilde{S}}$  be any Markov kernel on  $\mathcal{W}$  with access to the supersample  $\tilde{S}$  but not to the indices  $U$ . For all  $\lambda \in \mathbb{R}_+$  and every  $\beta \in (0, 1)$ , with probability no less than  $1 - \beta$*

$$\mathbb{E}^{\tilde{S}, U} \widehat{\mathcal{R}}(W, S_{\text{ghost}}) \leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \frac{1}{\lambda} \left[ D(\mathbb{P}_W^S \| \mathbb{Q}_W^{\tilde{S}}) + \log \frac{1}{\beta} \right] + \frac{\lambda c^2}{2n} \quad (20)$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

Comparing (8) and (20), it is apparent that the technique from Theorem 12 can be readily applied. However, as the function  $\psi$  is explicit here, to highlight the difference between this approach and the one quantizing the parameter space (Alquier, 2021, Section 2.1.4), we show how to proceed below. This also showcases when one would want to choose a more stringent event as anticipated in Section 3.2.2.

**Theorem 31** *Consider a bounded loss function  $\ell : \mathcal{W} \times \mathcal{Z} \rightarrow [a, b]$ . Let  $\mathbb{Q}_W^{\tilde{S}}$  be any Markov kernel on  $\mathcal{W}$  with access to the supersample  $\tilde{S}$  but not to the indices  $U$ . Then, for every  $\beta \in (0, 1)$ , with probability no smaller than  $1 - \beta$*

$$\mathbb{E}^S \mathcal{R}(W) \leq \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \sqrt{\frac{2(b-a)^2 (D(\mathbb{P}_W^S \| \mathbb{Q}_W^{\tilde{S}}) + \log \frac{en}{\beta})}{n}} + \sqrt{\frac{(b-a)^2 \log(\frac{4}{\beta})}{2n}}$$

holds simultaneously for every posterior  $\mathbb{P}_W^S$ .

**Proof** Similarly to the proof of Theorem 12, let  $\mathcal{B}_\lambda$  be the complement of the event in (20), and define the event  $\mathcal{E} := \{D(\mathbb{P}_W^S \| \mathbb{Q}_W^{\tilde{S}}) \leq k_{\max}\}$ , and the sub-events  $\mathcal{E}_1 := \{D(\mathbb{P}_W^S \| \mathbb{Q}_W^{\tilde{S}}) \leq 1\}$  and  $\mathcal{E}_k := \{[D(\mathbb{P}_W^S \| \mathbb{Q}_W^{\tilde{S}})] = k\}$  for  $k = 2, \dots, k_{\max}$ . Define also  $\mathcal{K} := \{k \in \mathbb{N} : 1 \leq k \leq k_{\max} \text{ and } \mathbb{P}[\mathcal{E}_k] > 0\}$ . Then, for all  $k \in \mathcal{K}$ , given the event  $\mathcal{E}_k$ , with probability at most  $\mathbb{P}[\mathcal{B}_\lambda | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^{\tilde{S}, U} \widehat{\mathcal{R}}(W, S_{\text{ghost}}) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \frac{1}{\lambda} \left[ k + \log \frac{1}{\beta} \right] + \frac{\lambda(b-a)^2}{2n}$$

for all  $\lambda > 0$ . The above equation is optimized for

$$\lambda = \lambda_k := \sqrt{\frac{2n}{(b-a)^2} \left( k + \log \frac{1}{\beta} \right)}$$

and using that  $k \leq D(\mathbb{P}_W^S \| \mathbb{Q}_W^{\tilde{S}}) + 1$  yields that, given the event  $\mathcal{E}_k$ , with probability smaller or equal than  $\mathbb{P}[\mathcal{B}_{\lambda_k} | \mathcal{E}_k]$ , there exists some posterior  $\mathbb{P}_W^S$  such that

$$\mathbb{E}^{\tilde{S}, U} \widehat{\mathcal{R}}(W, S_{\text{ghost}}) > \mathbb{E}^S \widehat{\mathcal{R}}(W, S) + \sqrt{\frac{2(b-a)^2 (D(\mathbb{P}_W^S \| \mathbb{Q}_W^{\tilde{S}}) + \log \frac{1}{\beta} + 1)}{n}}. \quad (21)$$

Let  $\mathcal{B}'$  be the event in (21) and note that  $\mathbb{P}[B' \cap E^c] = 0$  as long as  $k_{\max} \geq \frac{n}{2} - 1$  by the boundedness of the loss. Solving as in Theorem 12 and using (Hellström and Durisi, 2020, Theorem 3) completes the proof. ■

The parameter  $\lambda$  optimization strategies from Langford and Caruana (2001) and Catoni (2003) are based on a quantization of the *parameter space*. They choose a countable set  $\mathcal{A} \subset \mathbb{R}_+$  and use other techniques (e.g., rounding) to ensure that the optimization can be done on a larger set  $\mathcal{A}'$  with a union bound price of  $\log |\mathcal{A}|$  (Alquier, 2021, Section 2.1.4). However, it is not always certain that the optimal parameter lies on the extended set  $\mathcal{A}'$ , and thus a parameter-free closed form expression of the bound is unattainable. Here, we instead partition *the set of events described by the random variables*, create a further upper bound that depends only on deterministic terms for each event, and select a deterministic parameter for that event. Then, we pay a union bound price for each event considered. Therefore, as mentioned previously, our technique can be seen as an optimization of the set  $\mathcal{A}^* := \{\lambda_0, \dots, \lambda_{k_{\max}}\}$  of parameters for which we know that an *almost optimal* bound can be reached, where the distance to optimality is given by how loose is the upper bound depending only on deterministic terms (B.8).

## References

- Pierre Alquier. Transductive and inductive adaptative inference for regression and density estimation. *University Paris 6*, 2006.
- Pierre Alquier. User-friendly introduction to PAC-Bayes bounds. *arXiv preprint arXiv:2110.11216*, 2021.
- Pierre Alquier and Gérard Biau. Sparse single-index model. *Journal of Machine Learning Research*, 14(1), 2013.
- Pierre Alquier and Benjamin Guedj. Simpler PAC-Bayesian bounds for hostile data. *Machine Learning*, 107(5):887–902, 2018.
- Amiran Ambroladze, Emilio Parrado-Hernández, and John Shawe-Taylor. Tighter PAC-Bayes bounds. *Advances in neural information processing systems (NeurIPS)*, 19, 2006.
- Ron Amit, Baruch Epstein, Shay Moran, and Ron Meir. Integral probability metrics PAC-Bayes bounds. *Advances in neural information processing systems (NeurIPS)*, 35:3123–3136, 2022.
- Gautier Appert and Olivier Catoni. New bounds for  $k$ -means and information  $k$ -means. *arXiv preprint arXiv:2101.05728*, 2021.
- Søren Asmussen, Jens Ledet Jensen, and Leonardo Rojas-Nandayapa. On the Laplace transform of the lognormal distribution. *Methodology and Computing in Applied Probability*, 18:441–458, 2016.

- Pradeep Kr Banerjee and Guido Montúfar. Information complexity and generalization bounds. In *IEEE International Symposium on Information Theory (ISIT)*, pages 676–681. IEEE, 2021.
- Raef Bassily, Shay Moran, Ido Nachum, Jonathan Shafer, and Amir Yehudayoff. Learners that use little information. In *International Conference on Algorithmic Learning Theory (ALT)*, pages 25–55. PMLR, 2018.
- Luc Bégin, Pascal Germain, François Laviolette, and Jean-François Roy. PAC-Bayesian theory for transductive learning. In *Artificial Intelligence and Statistics*, pages 105–113. PMLR, 2014.
- Luc Bégin, Pascal Germain, François Laviolette, and Jean-François Roy. PAC-Bayesian bounds based on the Rényi divergence. In *Artificial Intelligence and Statistics*, pages 435–444. PMLR, 2016.
- Bernard Bercu and Abderrahmen Touati. Exponential inequalities for self-normalized martingales with applications. *Annals of Applied Probability*, 18:1848–1869, 2008.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- Yuheng Bu, Shaofeng Zou, and Venugopal V Veeravalli. Tightening mutual information-based bounds on generalization error. *IEEE Journal on Selected Areas in Information Theory*, 1(1):121–130, 2020.
- Olivier Catoni. A PAC-Bayesian approach to adaptive classification. *preprint*, 840, 2003.
- Olivier Catoni. *Statistical learning theory and stochastic optimization: Ecole d’Eté de Probabilités de Saint-Flour, Summer School XXXI-2001*, volume 1851. Springer Science & Business Media, 2004.
- Olivier Catoni. PAC-Bayesian supervised classification: The thermodynamics of statistical learning. *IMS Lecture Notes Monograph Series*, 56:163pp, 2007.
- Olivier Catoni. PAC-Bayes bounds for supervised classification. *Measures of Complexity: Festschrift for Alexey Chervonenkis*, pages 287–302, 2015.
- Ioannis Chatzigeorgiou. Bounds on the Lambert function and their application to the outage analysis of user cooperation. *IEEE Communications Letters*, 17(8):1505–1508, 2013.
- Badr-Eddine Chérif-Abdellatif, Yuyang Shi, Arnaud Doucet, and Benjamin Guedj. On PAC-Bayesian reconstruction guarantees for VAEs. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 3066–3079. PMLR, 2022.
- Ben Chugg, Hongjian Wang, and Aaditya Ramdas. A unified recipe for deriving (time-uniform) pac-bayes bounds. *Journal of Machine Learning Research*, 24(372):1–61, 2023.
- DA Darling and Herbert Robbins. Iterated logarithm inequalities. *Proceedings of the National Academy of Sciences*, 57(5):1188–1192, 1967.

- Victor H de la Peña, Michael J Klass Lai, and Tze Leung Lai. Pseudo-maximization and self-normalized processes. *Probability Surveys*, 4:172–192, 2007.
- Monroe D Donsker and SR Srinivasa Varadhan. Asymptotic evaluation of certain Markov process expectations for large time, I. *Communications on Pure and Applied Mathematics*, 28(1):1–47, 1975.
- Gintare Karolina Dziugaite and Daniel M. Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2017.
- Gintare Karolina Dziugaite and Daniel M. Roy. Data-dependent PAC-Bayes priors via differential privacy. *Advances in neural information processing systems (NeurIPS)*, 31, 2018.
- Amedeo Roberto Esposito, Michael Gastpar, and Ibrahim Issa. Generalization error bounds via Rényi-, f-divergences and maximal leakage. *IEEE Transactions on Information Theory*, 67(8):4986–5004, 2021.
- Andrew Foong, Wessel Bruinsma, David Burt, and Richard Turner. How tight can PAC-Bayes be in the small data regime? *Advances in neural information processing systems (NeurIPS)*, 34:4093–4105, 2021.
- Pascal Germain, Alexandre Lacasse, François Laviolette, and Mario Marchand. PAC-Bayesian learning of linear classifiers. In *International Conference on Machine Learning (ICML)*, pages 353–360, 2009.
- Pascal Germain, Alexandre Lacasse, Francois Laviolette, Mario March, and Jean-Francois Roy. Risk bounds for the majority vote: From a PAC-Bayesian analysis to a learning algorithm. *Journal of Machine Learning Research*, 16(26):787–860, 2015.
- Josiah Willard Gibbs. *Elementary principles in statistical mechanics: developed with especial reference to the rational foundations of thermodynamics*. C. Scribner’s sons, 1902.
- Peter Grunwald, Thomas Steinke, and Lydia Zakyntinou. PAC-Bayes, MAC-Bayes and conditional mutual information: Fast rate bounds that handle general VC classes. In *Conference on Learning Theory (COLT)*, pages 2217–2247. PMLR, 2021.
- Benjamin Guedj and Pierre Alquier. PAC-Bayesian estimation and prediction in sparse additive models. *Electronic Journal of Statistics*, 7:264–291, 2013.
- Benjamin Guedj and Louis Pujol. Still no free lunches: the price to pay for tighter PAC-Bayes bounds. *Entropy*, 23(11):1529, 2021.
- Maxime Haddouche and Benjamin Guedj. PAC-Bayes generalisation bounds for heavy-tailed losses through supermartingales. *Transactions on Machine Learning Research*, 2023a. ISSN 2835-8856. URL <https://openreview.net/forum?id=qxrwt6F3sf>.
- Maxime Haddouche and Benjamin Guedj. Wasserstein PAC-Bayes learning: A bridge between generalisation and optimisation. *arXiv preprint arXiv:2304.07048*, 2023b.

- Maxime Haddouche, Benjamin Guedj, Omar Rivasplata, and John Shawe-Taylor. PAC-Bayes unleashed: Generalisation bounds with unbounded losses. *Entropy*, 23(10):1330, 2021.
- Mahdi Haghifam, Gintare Karolina Dziugaite, Shay Moran, and Dan Roy. Towards a unified information-theoretic framework for generalization. *Advances in neural information processing systems (NeurIPS)*, 34:26370–26381, 2021.
- Mahdi Haghifam, Borja Rodríguez-Gálvez, Ragnar Thobaben, Mikael Skoglund, Daniel M Roy, and Gintare Karolina Dziugaite. Limitations of information-theoretic generalization bounds for gradient descent methods in stochastic convex optimization. In *International Conference on Algorithmic Learning Theory (ALT)*, pages 663–706. PMLR, 2023.
- Fredrik Hellström and Giuseppe Durisi. Generalization bounds via information density and conditional information density. *IEEE Journal on Selected Areas in Information Theory*, 1(3):824–839, 2020.
- Fredrik Hellström and Giuseppe Durisi. Corrections to “Generalization bounds via information density and conditional information density”. *IEEE Journal on Selected Areas in Information Theory*, 2(3):1072–1073, 2021.
- Fredrik Hellström and Giuseppe Durisi. A new family of generalization bounds using sample-wise evaluated CMI. In *Advances in neural information processing systems (NeurIPS)*, 2022.
- Matthew Higgs and John Shawe-Taylor. A PAC-Bayes bound for tailored density estimation. In *International Conference on Algorithmic Learning Theory (ALT)*, pages 148–162. Springer, 2010.
- Matthew Holland. PAC-Bayes under potentially heavy tails. *Advances in neural information processing systems (NeurIPS)*, 32, 2019.
- Steven R Howard, Aaditya Ramdas, Jon McAuliffe, and Jasjeet Sekhon. Time-uniform, nonparametric, nonasymptotic confidence sequences. *The Annals of Statistics*, 49(2):1055–1080, 2021.
- Kyoungseok Jang, Kwang-Sung Jun, Ilja Kuzborskij, and Francesco Orabona. Tighter pac-bayes bounds through coin-betting. In Gergely Neu and Lorenzo Rosasco, editors, *Conference on Learning Theory (COLT)*, volume 195 of *Proceedings of Machine Learning Research*, pages 2240–2264. PMLR, 12–15 Jul 2023. URL <https://proceedings.mlr.press/v195/jang23a.html>.
- Sham M Kakade, Karthik Sridharan, and Ambuj Tewari. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. *Advances in neural information processing systems (NeurIPS)*, 21, 2008.
- Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best arm identification in multi-armed bandit models. *Journal of Machine Learning Research*, 17:1–42, 2016.



- Ilja Kuzborskij and Csaba Szepesvári. Efron-Stein PAC-Bayesian inequalities. *arXiv preprint arXiv:1909.01931*, 2019.
- Ilja Kuzborskij, Kwang-Sung Jun, Yulian Wu, Kyoungseok Jang, and Francesco Orabona. Better-than-KL PAC-Bayes bounds. *arXiv preprint arXiv:2402.09201*, 2024.
- John Langford and Rich Caruana. (Not) bounding the true error. *Advances in neural information processing systems (NeurIPS)*, 14, 2001.
- John Langford and Matthias Seeger. Bounds for averaging classifiers. Technical report, School of Computer Science, Carnegie Mellon University, 2001.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Roi Livni and Shay Moran. A limitation of the PAC-Bayes framework. *Advances in neural information processing systems (NeurIPS)*, 33:20543–20553, 2020.
- Sanae Lotfi, Marc Finzi, Sanyam Kapoor, Andres Potapczynski, Micah Goldblum, and Andrew G Wilson. PAC-Bayes compression bounds so tight that they can explain generalization. *Advances in neural information processing systems (NeurIPS)*, 35:31459–31473, 2022.
- The Tien Mai and Pierre Alquier. A Bayesian approach for noisy matrix completion: Optimal rate under general sampling distribution. *Electronic Journal of Statistics*, 9(1):823–841, 2015.
- Katalin Marton. A measure concentration inequality for contracting Markov chains. *Geometric & Functional Analysis GFA*, 6(3):556–571, 1996.
- Andreas Maurer. A note on the PAC Bayesian theorem. *arXiv preprint cs/0411099*, 2004.
- David A McAllester. Some PAC-Bayesian theorems. In *Conference on Computational learning theory (COLT)*, pages 230–234, 1998.
- David A McAllester. PAC-Bayesian model averaging. In *Conference on Computational learning theory (COLT)*, pages 164–170, 1999.
- David A McAllester. PAC-Bayesian stochastic model selection. *Machine Learning*, 51(1): 5–21, 2003.
- Zakaria Mhammedi, Peter Grünwald, and Benjamin Guedj. PAC-Bayes un-expected Bernstein inequality. *Advances in neural information processing systems (NeurIPS)*, 32, 2019.
- Ido Nachum, Jonathan Shafer, Thomas Weinberger, and Michael Gastpar. Fantastic generalization measures are nowhere to be found. In *International Conference on Learning Representations (ICLR)*, 2023.
- Vaishnavh Nagarajan and J Zico Kolter. Uniform convergence may be unable to explain generalization in deep learning. *Advances in neural information processing systems (NeurIPS)*, 32, 2019.

- Balas K Natarajan. On learning sets and functions. *Machine Learning*, 4(1):67–97, 1989.
- N. Balakrishnan Norman L. Johnson, Samuel Kotz. *Continuous Univariate Distributions*, volume 1. John Wiley & Sons Inc., 1994.
- Kento Nozawa and Issei Sato. PAC-Bayes analysis of sentence representation. *arXiv preprint arXiv:1902.04247*, 2019.
- Kento Nozawa, Pascal Germain, and Benjamin Guedj. PAC-Bayesian contrastive unsupervised representation learning. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 21–30. PMLR, 2020.
- Yuki Ohnishi and Jean Honorio. Novel change of measure inequalities with applications to PAC-Bayesian bounds and Monte Carlo estimation. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 1711–1719. PMLR, 2021.
- María Pérez-Ortiz, Omar Rivasplata, John Shawe-Taylor, and Csaba Szepesvári. Tighter risk certificates for neural networks. *The Journal of Machine Learning Research*, 22(1):10326–10365, 2021.
- Yury Polyanskiy and Yihong Wu. *Information Theory: From Coding to Learning*. Cambridge University Press, 1st edition, 2023.
- Liva Ralaivola, Marie Szafranski, and Guillaume Stempfel. Chromatic PAC-Bayes bounds for non-iid data: Applications to ranking and stationary  $\beta$ -mixing processes. *The Journal of Machine Learning Research*, 11:1927–1956, 2010.
- David Reeb, Andreas Doerr, Sebastian Gerwinn, and Barbara Rakitsch. Learning Gaussian processes by minimizing PAC-Bayesian generalization bounds. *Advances in neural information processing systems (NeurIPS)*, 31, 2018.
- Omar Rivasplata, Vikram M Tankasali, and Csaba Szepesvári. PAC-Bayes with backprop. *arXiv preprint arXiv:1908.07380*, 2019.
- Herbert Robbins and David Siegmund. Iterated logarithm inequalities and related statistical procedures. *Mathematics of the Decision Sciences*, 2:267–279, 1968.
- Matthias Seeger. PAC-Bayesian generalisation error bounds for Gaussian process classification. *Journal of machine learning research*, 3(Oct):233–269, 2002.
- Yevgeny Seldin. Machine learning. the science of selection under uncertainty. *Lecture Notes*, 2023. URL <https://sites.google.com/site/yevgenyseldin/teaching?authuser=0>.
- Yevgeny Seldin and Naftali Tishby. PAC-Bayesian analysis of co-clustering and beyond. *Journal of Machine Learning Research*, 11(12), 2010.
- Yevgeny Seldin, François Laviolette, Nicolo Cesa-Bianchi, John Shawe-Taylor, and Peter Auer. PAC-Bayesian inequalities for martingales. *IEEE Transactions on Information Theory*, 58(12):7086–7093, 2012.

- Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- John Shawe-Taylor, Peter L Bartlett, Robert C Williamson, and Martin Anthony. A framework for structural risk minimisation. In *Conference on Computational learning theory (COLT)*, pages 68–76, 1996.
- Thomas Steinke and Lydia Zakynthinou. Reasoning about generalization via conditional mutual information. In *Conference on Learning Theory (COLT)*, pages 3437–3452. PMLR, 2020.
- Niklas Thiemann, Christian Igel, Olivier Wintenberger, and Yevgeny Seldin. A strongly quasiconvex PAC-Bayesian bound. In *International Conference on Algorithmic Learning Theory (ALT)*, pages 466–492. PMLR, 2017.
- Ilya O Tolstikhin and Yevgeny Seldin. PAC-Bayes-empirical-Bernstein inequality. *Advances in neural information processing systems (NeurIPS)*, 26, 2013.
- Paul Viallard, Maxime Haddouche, Umut Simsekli, and Benjamin Guedj. Learning via Wasserstein-based high probability generalisation bounds. *Advances in neural information processing systems (NeurIPS)*, 36, 2024a.
- Paul Viallard, Maxime Haddouche, Umut Şimşekli, and Benjamin Guedj. Tighter generalisation bounds via interpolation. *arXiv preprint arXiv:2402.05101*, 2024b.
- Jean Ville. Étude critique de la notion de collectif. *Bull. Amer. Math. Soc.*, 45(11):824, 1939.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- Zhen Wang, Luming Shen, Yu Miao, Shanshan Chen, and Wenfei Xu. PAC-Bayesian inequalities of some random variables sequences. *Journal of Inequalities and Applications*, 2015(1):1–8, 2015.
- Yi-Shan Wu and Yevgeny Seldin. Split-kl and PAC-Bayes-split-kl inequalities for ternary random variables. *Advances in neural information processing systems (NeurIPS)*, 35: 11369–11381, 2022.
- Yi-Shan Wu, Andres Masegosa, Stephan Lorenzen, Christian Igel, and Yevgeny Seldin. Chebyshev-Cantelli PAC-Bayes-Bennett inequality for the weighted majority vote. *Advances in neural information processing systems (NeurIPS)*, 34:12625–12636, 2021.
- Tong Zhang. Information-theoretic upper and lower bounds for statistical estimation. *IEEE Transactions on Information Theory*, 52(4):1307–1321, 2006.
- Wenda Zhou, Victor Veitch, Morgane Austern, Ryan P Adams, and Peter Orbanz. Non-vacuous generalization bounds at the im-agenet scale: A PAC-Bayesian compression approach. In *International Conference on Learning Representations (ICLR)*, 2019.