



Ph.D. Thesis

Yi-Shan Wu

Second-Order Concentration Inequalities with Application to the Weighted Majority Vote

Advisors: Yevgeny Seldin and Anders Søgaard

This thesis has been submitted to the Ph.D. School of The Faculty of Science, University of Copenhagen on March 15th 2023.

Abstract

The weighted majority vote is part of the winning strategies in many machine learning competitions. It is an integral part of random forests and boosting and is also used to combine predictions of heterogeneous classifiers. While generalization guarantees and theoretically grounded optimization algorithms of a weighted majority vote have been a research topic for decades, it is still unclear how to achieve a tight generalization bound that also serves as a faithful optimization objective for achieving good test performance of the resulting weighted majority vote.

There have been extensive studies into second-order analysis that take into account the correlation of errors of the classifiers, which is the key power of the weighted majority vote. However, the existing bounds are either tight but hard to optimize, e.g., the C-bounds, or easy to optimize but come at the cost of being less tight, e.g., the tandem bound. In this thesis, we derive new second-order oracle bounds, where we show it's possible to achieve bounds that are as tight as the C-bounds but remain easy to optimize as the tandem bound.

The concentration of measure inequalities are also needed to transform oracle bounds into empirical bounds. In particular, we derive a new second-order concentration inequality for ternary random variables, which are random variables with a support size of three. Ternary random variables appear in various applications, including the analysis of the weighted majority vote, excess losses, and learning with abstention. The new inequality takes advantage of both the kl inequality that exploits the combinatorial structure in the case of binary random variables, and the Bernstein-type inequalities that are effective if the random variables have low variance.

In this thesis, our focus is on ensemble classifiers and randomized classifiers. To handle this, we use PAC-Bayesian analysis to convert oracle bounds to empirical bounds. This involves applying standard PAC-Bayesian techniques to obtain the PAC-Bayesian versions of various concentration inequalities for random variables.

Resumé

Den weighted majority vote er en del af de vindende strategier i mange konkurrencer inden for machine learning. Det er en integreret del af random forests og boosting og bruges også til at kombinere forudsigelser fra heterogene klassifikatorer. Mens generaliseringsgarantier og teoretisk funderede optimeringsalgoritmer af en weighted majority vote har været et forskningsemne i årtier, er det stadig uklart, hvordan man opnår en stram generaliseringsgrænse, der også fungerer som en pålidelig optimeringsmål for at opnå god testydelse af det resulterende weighted majority vote.

Der er blevet omfattende undersøgelser af andenordensanalyse, der tager hensyn til korrelationen af fejl fra klassifikatorerne, som er den centrale styrke i den weighted majority vote. Imidlertid er eksisterende grænser enten stramme, men svære at optimere, f.eks. C-bounds, eller lette at optimere, men på bekostning af at være mindre stramme, f.eks. tandem bound. I denne afhandling udleder vi nye andenordens orakelgrænser, hvor vi viser, at det er muligt at opnå grænser, der er lige så stramme som C-bounds, men stadig er lette at optimere som tandem bound.

Koncentration uligheder er også nødvendige for at omdanne orakelgrænser til empiriske grænser. Vi udleder især en ny andenordens koncentration uligheder for ternære stokastiske variable, som er stokastiske variable med tre mulige udfald. Ternære stokastiske variable optræder i forskellige anvendelser, herunder analysen af den weighted majority vote, excess losses, og learning with abstention. Den nye ulighed udnytter både kl-uligheden, som udnytter den kombinatoriske struktur i stokastiske variable, og de uligheder af Bernstein-typen, der er effektive, hvis de stokastiske variable har lav varians.

I denne afhandling fokuserer vi på ensembleklassifikatorer og randomiserede klassifikatorer. For at håndtere dette bruger vi PAC-Bayesian analyse til at konvertere orakelgrænser til empiriske grænser. Dette indebærer anvendelse af standard PAC-Bayesian teknikker til at opnå PAC-Bayesian versioner af forskellige koncentration uligheder for stokastiske variable.

Acknowledgements

I would like to start by expressing my deep gratitude to my supervisors, Yevgeny Seldin and Anders Søgaard, and in particular, Yevgeny Seldin, who has played a vital role in my academic journey. Yevgeny has consistently provided me with support and encouragement in both my academic pursuits and personal life in Denmark, and it has been an honor and privilege to be able to work with him. I would also like to express my appreciation to Andres R. Masegosa for our close collaborations, which have been a great source of inspiration.

I am fortunate to have had the opportunity to visit Peter Grünwald and Wouter M. Koolen at CWI on several occasions. I extend my sincere thanks to Peter and Wouter for their kind hospitality and the engaging discussions we have had. I would also like to thank Omar Rivasplata for his warm welcome and introduction to various research topics during my visit to London.

Moreover, I am grateful for the wonderful colleagues I have had the pleasure of being with during my time at DIKU, especially the Ph.D. fellows from the DeLTA group, Chloé, Saeed, Yunlian, Yijie, and Hippolyte. Their presence and support have been invaluable. I am also thankful for my Taiwanese friends in Europe, who have been a great help in adjusting to life in Europe.

Additionally, I want to take a moment to recognize my previous colleagues at IIS, Academia Sinica in Taiwan, particularly my former advisor Chi-Jen Lu, who guided me into the fascinating world of machine learning research and provided me with patient guidance and assistance. I would also like to thank my academic partner, Chen Yanlin, for always engaging in insightful discussions. Additionally, I want to acknowledge Wang, Po-An, and Lin, Jing-Hua for introducing me to the delightful world of coffee-making and coffee-tasting, which I have since shared with others.

Finally, I would like to express my deepest appreciation to my family and friends, especially my mother, sister, and aunt, for their unwavering support. And to my father, I know you are always there with me.

Table of Contents

Abstract	ii
Resumé	iii
Acknowledgements	iv
1 Introduction	1
1.1 Outline of the Thesis	2
1.2 Main Contributions	3
2 Chebyshev-Cantelli PAC-Bayes-Bennett Inequality for the Weighted Majority Vote	5
2.1 Introduction	6
2.2 Problem setup	8
2.3 A review of prior first and second-order oracle bounds	9
2.4 Main Contributions	10
2.5 From oracle to empirical bounds	14
2.6 Experiments	17
2.7 Discussion	21
2.8 Appendix	21
3 Split-kl and PAC-Bayes-split-kl Inequalities for Ternary Random Variables	38
3.1 Introduction	39
3.2 Concentration of Measure Inequalities for Sums of Independent Random Variables	41
3.3 PAC-Bayesian Inequalities	45
3.4 Experiments	51
3.5 Discussion	53
3.6 Appendix	54

4 Summary and Discussion	77
List of Publications	79
Bibliography	80

Chapter 1

Introduction

Ensemble methods are machine learning methods that combine the predictions of multiple classifiers to produce more accurate and robust models. One such method is the weighted majority vote, which assigns a weight to each base classifier, or voter, and outputs the majority of the predictions based on the assigned weights. It is an integral part of random forests and boosting, and is also applied to combining the predictions of heterogeneous classifiers. The power of the majority vote is in the cancellation of errors effect: when the errors of individual classifiers are independent or anticorrelated and the error probability of individual classifiers is smaller, then the errors average out and the majority vote tends to outperform the individual classifiers.

While generalization guarantees of a weighted majority vote have been a research topic for decades, with optimal weighting derived under certain assumptions (Berend and Kontorovich, 2016), the assumptions are typically not satisfied in practice. Therefore, techniques that estimate and solve for the optimal combination of voter predictions from the available data are necessary. The most basic result is the first-order oracle bound, which bounds the expected loss of weighted majority vote by the average of expected losses of the individual classifiers. However, due to ignoring error correlation among classifiers, the method assigns excessive weight to the best-performing ones, reducing the majority vote to just a few classifiers and significantly increasing the test error (Lorenzen et al., 2019).

To take the correlation of errors into consideration, previous works proposed second-order oracle bounds, the C-bounds (Lacasse et al., 2007; Germain et al., 2015; Laviolette et al., 2017), and the tandem bound (Masegosa et al., 2020). While the

C-bounds are tight, they are hard to estimate and optimize. On the other hand, the tandem bound is easy to optimize but is not as tight as the C-bounds.

One of the focuses of this thesis is to bridge between these two second-order oracle bounds, where we show it's possible to achieve bounds that are as tight as the C-bounds but easy to optimize as the tandem bound.

Since accessing oracle bounds is often impossible, we have to use concentration inequalities to transform them into empirical bounds. In particular, an example of the oracle quantities that requires concentration inequalities is a ternary random variable, which is a random variable with a support size of three. It appears in various applications, including the analysis of the weighted majority vote, excess losses, and learning with abstention. Ternary random variables are slightly more complex than binary random variables but simpler than general bounded random variables.

While the kl inequality (Maurer, 2004; Langford, 2005) is known to be tight for binary random variables and applicable to any bounded random variables, it is not necessarily a good choice for bounded random variables that can take more than two values. Instead, the Bernstein-type inequalities (Maurer and Pontil, 2009; Cesa-Bianchi et al., 2007; Mhammedi et al., 2019) can be more effective due to their ability to exploit low variance. However, they can be loose for binary random variables. Previous studies (Tolstikhin and Seldin, 2013; Mhammedi et al., 2019) have attempted to balance the need for exploiting low variance while still achieving the same level of tightness as the kl inequality if a distribution happens to be close to binary, but the problem remains an open question.

Another focus of this thesis is to resolve this question for the case of ternary random variables by a new concentration inequality based on the kl inequality.

Finally, because we are dealing with ensemble classifiers and randomized classifiers, we utilize PAC-Bayesian analysis to transform the oracle bounds into empirical bounds. This involves using a standard PAC-Bayesian technique to derive the PAC-Bayesian versions of many concentration inequalities for random variables.

1.1 Outline of the Thesis

The following two chapters contain the papers Wu et al. (2021) and Wu and Seldin (2022).

Chapter 2 corresponds to Wu et al. (2021), where we present a new second-order

oracle bound for the expected risk of a weighted majority vote. The bound is based on a novel parametric form of the Chebyshev-Cantelli inequality (a.k.a. one-sided Chebyshev's), which is amenable to efficient minimization. The new form resolves the optimization challenge faced by prior oracle bounds based on the Chebyshev-Cantelli inequality, the C-bounds (Germain et al., 2015), and, at the same time, it improves on the oracle bound based on second-order Markov's inequality introduced by Masegosa et al. (2020). We also derive a new second-order concentration of measure inequality, which we name PAC-Bayes-Bennett, since it combines PAC-Bayesian bounding with Bennett's inequality. We use it for empirical estimation of the oracle bound. The PAC-Bayes-Bennett inequality improves on the PAC-Bayes-Bernstein inequality of Seldin et al. (2012). We provide an empirical evaluation demonstrating that the new bounds can improve on the previous second-order bounds for weighted majority vote.

Chapter 3 corresponds to Wu and Seldin (2022), where we present a new concentration of measure inequality for sums of independent bounded random variables, which we name a split-kl inequality. The inequality is particularly well-suited for ternary random variables, which appear in various problems, including analysis of weighted majority votes, analysis of excess losses in classification, and learning with abstention. We demonstrate that for ternary random variables the inequality is simultaneously competitive with the kl inequality, the Empirical Bernstein inequality, and the Unexpected Bernstein inequality, and in certain regimes outperforms all of them. It resolves an open question by Tolstikhin and Seldin (2013) and Mhammedi et al. (2019) on how to match simultaneously the combinatorial power of the kl inequality when the distribution happens to be close to binary and the power of Bernstein inequalities to exploit low variance when the probability mass is concentrated on the middle value. We also derive a PAC-Bayes-split-kl inequality and compare it with the PAC-Bayes-kl, PAC-Bayes-Empirical-Bennett, and PAC-Bayes-Unexpected-Bernstein inequalities in an analysis of excess losses and in an analysis of a weighted majority vote for several UCI datasets. Last but not least, our study provides the first direct comparison of the Empirical Bernstein and Unexpected Bernstein inequalities and their PAC-Bayes extensions.

We finish this work by a discussion of these results in Chapter 4.

1.2 Main Contributions

The main contributions of this work are:

- We propose a new parametric form of the Chebyshev-Cantelli inequality, which

has no variance in the denominator and preserves tightness of the original bound. The new form allows efficient minimization and empirical estimation.

- We propose two new second-order oracle bounds for the weighted majority vote based on the new form of the Chebyshev-Cantelli inequality. The bounds have two advantages: (1) they are amenable to tight translation to empirical bounds; and (2) the resulting empirical bounds are amenable to efficient minimization.
- We propose four new second-order empirical bounds for the weighted majority vote based on the two newly proposed second-order oracle bounds. The empirical bounds are amenable to efficient minimization. We show they lead to tighter bounds and better test performance than the existing second-order bounds.
- We derive a new concentration of measure inequality, which we name the PAC-Bayes-Bennett inequality. It can be applied to arbitrary loss functions taking values in an interval with bounded length. Also, it improves on the PAC-Bayes-Bernstein inequality of Seldin et al. (2012).
- We derive a new concentration of measure inequality, which we name the split-kl inequality. We also derive its PAC-Bayes counterpart, the PAC-Bayes-split-kl inequality. They can be applied to arbitrary bounded loss functions and are particularly well-suited for ternary random variables.
- Importantly, we demonstrate that for ternary random variables, the split-kl inequality is simultaneously competitive with the kl inequality, the Empirical Bernstein inequality, and the Unexpected Bernstein inequality, and in certain regimes outperforms all of them. It resolves an open question by Tolstikhin and Seldin (2013) and Mhammedi et al. (2019) on how to match simultaneously the combinatorial power of the kl inequality when the distribution happens to be close to binary and the power of Bernstein inequalities to exploit low variance when the probability mass is concentrated on the middle value.
- To the best of our knowledge, this is the first time when the Empirical Bernstein and the Unexpected Bernstein inequalities are directly compared, with and without the PAC-Bayesian extension. We also show that an inequality introduced by Cesa-Bianchi et al. (2007) yields a relaxation of the Unexpected Bernstein inequality by Mhammedi et al. (2019).
- We use the ideas of excess losses and informed priors (Ambroladze et al., 2007; Mhammedi et al., 2019), together with the proposed PAC-Bayes inequalities to improve the performance of the binary classification problems with linear classifiers considered by Mhammedi et al. (2019).

Chapter 2

Chebyshev-Cantelli PAC-Bayes-Bennett Inequality for the Weighted Majority Vote

The work presented in this chapter is based on a paper that has been published as:

Yi-Shan Wu, Andres Masegosa, Stephan Sloth Lorenzen, Christian Igel, and Yevgeny Seldin. Chebyshev-cantelli pac-bayes-bennett inequality for the weighted majority vote. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

Abstract

We present a new second-order oracle bound for the expected risk of a weighted majority vote. The bound is based on a novel parametric form of the Chebyshev-Cantelli inequality (a.k.a. one-sided Chebyshev's), which is amenable to efficient minimization. The new form resolves the optimization challenge faced by prior oracle bounds based on the Chebyshev-Cantelli inequality, the C-bounds (Germain et al., 2015), and, at the same time, it improves on the oracle bound based on second-order Markov's inequality introduced by Masegosa et al. (2020). We also derive a new concentration of measure inequality, which we name PAC-Bayes-Bennett, since it combines PAC-Bayesian bounding with Bennett's inequality. We use it for empirical estimation of the oracle bound. The PAC-Bayes-Bennett inequality improves on the PAC-Bayes-Bernstein inequality of Seldin et al. (2012). We provide an empirical evaluation demonstrating that the new bounds can improve on the work of Masegosa et al. (2020). Both the parametric form of the Chebyshev-Cantelli inequality and the PAC-Bayes-Bennett inequality may be of independent interest for the study of concentration of measure in other domains.

2.1 Introduction

Weighted majority vote is a central technique for combining predictions of multiple classifiers. It is an integral part of random forests (Breiman, 1996, 2001), boosting (Freund and Schapire, 1996), gradient boosting (Friedman, 1999, 2001; Mason et al., 1999; Chen and Guestrin, 2016), and it is also used to combine predictions of heterogeneous classifiers. It is part of the winning strategies in many machine learning competitions. The power of the majority vote is in the cancellation of errors effect: when the errors of individual classifiers are independent or anticorrelated and the error probability of individual classifiers is smaller than 0.5, then the errors average out and the majority vote tends to outperform the individual classifiers.

Generalization bounds for weighted majority vote and theoretically-grounded approaches to weight-tuning are decades-old research topics. Berend and Kontorovich (2016) derived an optimal solution under the assumption of known error rates and independence of errors of individual classifiers, but neither of the two assumptions is typically satisfied in practice.

In absence of the independence assumption, the most basic result is the first-order oracle bound, which is based on Markov's inequality and bounds the expected loss of

ρ -weighted majority vote by twice the ρ -weighted average of expected losses of the individual classifiers. This finding is so old and basic that Langford and Shawe-Taylor (2002) call it “the folk theorem”. The ρ -weighted average of the expected losses is then bounded using PAC-Bayesian bounds, turning the oracle bound into an empirical bound (McAllester, 1998; Seeger, 2002; Langford and Shawe-Taylor, 2002). While the translation from oracle to empirical bound is quite tight (Germain et al., 2009; Thiemann et al., 2017), the first-order oracle bound ignores the correlation of errors, which is the main power of the majority vote. As a result, its minimization overconcentrates the weights on the best-performing classifiers, effectively reducing the majority vote to very few or even a single best classifier, which leads to a significant deterioration of the test error (Lorenzen et al., 2019; Masegosa et al., 2020).

In order to take correlation of errors into account, Lacasse et al. (2007) derived second-order oracle bounds, the C-bounds, which are based on the Chebyshev-Cantelli inequality. The ideas were further developed by Laviolette et al. (2011), Germain et al. (2015), and Laviolette et al. (2017). However, they were only able to optimize the bounds in the highly restrictive setting of binary classification with self-complemented sets of classifiers and aligned priors and posteriors (Germain et al., 2015). Several follow-up works resorted to minimization of heuristic surrogates rather than the bound itself (Bauvin et al., 2020; Viallard et al., 2021). Furthermore, second-order oracle quantities in the denominator of the oracle bounds lead to looseness in their translation to empirical bounds (Lorenzen et al., 2019).

Masegosa et al. (2020) proposed an alternative second-order oracle bound, the tandem bound, based on second-order Markov’s inequality. While the second-order Markov’s inequality is weaker than the Chebyshev-Cantelli inequality, the resulting bound has no oracle quantities in the denominator, which allows tight translation to an empirical bound. Additionally, Masegosa et al. proposed an efficient procedure for minimization of their empirical bound. They have shown that minimization of the empirical bound does not lead to deterioration of the test error.

Our work can be seen as a bridge between the tandem bound and the C-bounds, and as an improvement of both. The key novelty is a new parametric form of Chebyshev-Cantelli inequality, which preserves the tightness of Chebyshev-Cantelli, but avoids oracle quantities in the denominator. This allows both efficient translation to empirical bounds and efficient minimization. We derive two new second-order oracle bounds based on the new inequality, one using the tandem loss and the other using the tandem loss with an offset. For empirical estimation of the latter we derive a PAC-Bayes-Bennett inequality. The overall contributions can be summarized as follows:

1. We propose a new parametric form of the Chebyshev-Cantelli inequality, which has no variance in the denominator and preserves tightness of the original bound. The new form allows efficient minimization and empirical estimation.
2. We propose two new second-order oracle bounds for the weighted majority vote based on the new form of the Chebyshev-Cantelli inequality. The bounds have two advantages: (1) they are amenable to tight translation to empirical bounds; and (2) the resulting empirical bounds are amenable to efficient minimization.
3. We derive a new concentration of measure inequality, which we name the PAC-Bayes-Bennett inequality. It improves on the PAC-Bayes-Bernstein inequality of Seldin et al. (2012). We use the inequality for bounding the tandem loss with an offset.

2.2 Problem setup

The problem setup and notations are borrowed from Masegosa et al. (2020).

Multiclass classification. We let $S = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ be an independent identically distributed sample from $\mathcal{X} \times \mathcal{Y}$, drawn according to an unknown distribution D , where \mathcal{Y} is finite and \mathcal{X} is arbitrary. A hypothesis is a function $h : \mathcal{X} \rightarrow \mathcal{Y}$, and \mathcal{H} denotes a space of hypotheses. We evaluate the quality of h using the zero-one loss $\ell(h(X), Y) = \mathbf{1}(h(X) \neq Y)$, where $\mathbf{1}(\cdot)$ is the indicator function. The expected loss of h is denoted by $L(h) = \mathbb{E}_{(X,Y) \sim D}[\ell(h(X), Y)]$ and the empirical loss of h on a sample S of size n is denoted by $\hat{L}(h, S) = \frac{1}{n} \sum_{i=1}^n \ell(h(X_i), Y_i)$.

Randomized classifiers. A *randomized classifier* (a.k.a. Gibbs classifier) associated with a distribution ρ on \mathcal{H} , for each input X randomly draws a hypothesis $h \in \mathcal{H}$ according to ρ and predicts $h(X)$. The expected loss of a randomized classifier is given by $\mathbb{E}_{h \sim \rho}[L(h)]$ and the empirical loss by $\mathbb{E}_{h \sim \rho}[\hat{L}(h, S)]$. To simplify the notation we use $\mathbb{E}_D[\cdot]$ as a shorthand for $\mathbb{E}_{(X,Y) \sim D}[\cdot]$ and $\mathbb{E}_\rho[\cdot]$ as a shorthand for $\mathbb{E}_{h \sim \rho}[\cdot]$.

Ensemble classifiers and majority vote. Ensemble classifiers predict by taking a weighted aggregation of predictions by hypotheses from \mathcal{H} . The ρ -weighted majority vote MV_ρ predicts $MV_\rho(X) = \arg \max_{y \in \mathcal{Y}} \mathbb{E}_\rho[\mathbf{1}(h(X) = y)]$, where ties can be resolved arbitrarily.

2.3 A review of prior first and second-order oracle bounds

If majority voting makes an error, we know that at least a ρ -weighted half of the classifiers have made an error and, therefore, $\ell(\text{MV}_\rho(X), Y) \leq \mathbb{1}(\mathbb{E}_\rho[\mathbb{1}(h(X) \neq Y)] \geq 0.5)$. This observation leads to the well-known first-order oracle bound for the loss of weighted majority vote.

Theorem 2.1 (First-Order Oracle Bound).

$$L(\text{MV}_\rho) \leq 2\mathbb{E}_\rho[L(h)].$$

Proof. We have $L(\text{MV}_\rho) = \mathbb{E}_D[\ell(\text{MV}_\rho(X), Y)] \leq \mathbb{P}(\mathbb{E}_\rho[\mathbb{1}(h(X) \neq Y)] \geq 0.5)$. By applying Markov's inequality to random variable $Z = \mathbb{E}_\rho[\mathbb{1}(h(X) \neq Y)]$ we have:

$$L(\text{MV}_\rho) \leq \mathbb{P}(\mathbb{E}_\rho[\mathbb{1}(h(X) \neq Y)] \geq 0.5) \leq 2\mathbb{E}_D[\mathbb{E}_\rho[\mathbb{1}(h(X) \neq Y)]] = 2\mathbb{E}_\rho[L(h)]. \quad \square$$

PAC-Bayesian analysis can be used to bound $\mathbb{E}_\rho[L(h)]$ in Theorem 2.1 in terms of $\mathbb{E}_\rho[\hat{L}(h, S)]$, thus turning the oracle bound into an empirical one. The disadvantage of the first-order approach is that $\mathbb{E}_\rho[L(h)]$ ignores correlations of predictions, which is the main power of the majority vote.

Masegosa et al. (2020) have used second-order Markov's inequality, by which for a non-negative random variable Z and $\varepsilon > 0$

$$\mathbb{P}(Z \geq \varepsilon) = \mathbb{P}(Z^2 \geq \varepsilon^2) \leq \frac{\mathbb{E}[Z^2]}{\varepsilon^2}.$$

For pairs of hypotheses h and h' they have defined the *tandem loss* $\ell(h(X), h'(X), Y) = \mathbb{1}(h(X) \neq Y \wedge h'(X) \neq Y) = \mathbb{1}(h(X) \neq Y)\mathbb{1}(h'(X) \neq Y)$, also termed *joint error* by Lacasse et al. (2007), which counts an error only if both h and h' err on a sample (X, Y) . The corresponding *expected tandem loss* is defined by

$$L(h, h') = \mathbb{E}_D[\mathbb{1}(h(X) \neq Y)\mathbb{1}(h'(X) \neq Y)].$$

Lacasse et al. (2007) and Masegosa et al. (2020) have shown that expectation of the second moment of the weighted loss equals expectation of the tandem loss. Using ρ^2 as a shorthand for the product distribution $\rho \times \rho$ over $\mathcal{H} \times \mathcal{H}$ and $\mathbb{E}_{\rho^2}[L(h, h')]$ as a shorthand for $\mathbb{E}_{h \sim \rho, h' \sim \rho}[L(h, h')]$, the result is the following.

Lemma 2.1 (Masegosa et al., 2020). *In multiclass classification*

$$\mathbb{E}_D[\mathbb{E}_\rho[\mathbf{1}(h(X) \neq Y)]^2] = \mathbb{E}_{\rho^2}[L(h, h')].$$

By combining second-order Markov's inequality with Lemma 2.1, Masegosa et al. have shown the following result.

Theorem 2.2 (Masegosa et al., 2020). *In multiclass classification*

$$L(\text{MV}_\rho) \leq 4\mathbb{E}_{\rho^2}[L(h, h')].$$

Lacasse et al. (2007) have used the Chebyshev-Cantelli inequality to derive a different form of a second-order oracle bound. We use $\mathbb{V}[Z]$ to denote the variance of a random variable Z in the statement of Chebyshev-Cantelli inequality.

Theorem 2.3 (Chebyshev-Cantelli inequality). *For $\varepsilon > 0$*

$$\mathbb{P}(Z - \mathbb{E}[Z] \geq \varepsilon) \leq \frac{\mathbb{V}[Z]}{\varepsilon^2 + \mathbb{V}[Z]}.$$

Theorem 2.3 together with Lemma 2.1 leads to the following result, known as the *C-bound*.

Theorem 2.4 (Lacasse et al., 2007; Masegosa et al., 2020). *If $\mathbb{E}_\rho[L(h)] \leq \frac{1}{2}$, then*

$$L(\text{MV}_\rho) \leq \frac{\mathbb{E}_{\rho^2}[L(h, h')] - \mathbb{E}_\rho[L(h)]^2}{\frac{1}{4} + \mathbb{E}_{\rho^2}[L(h, h')] - \mathbb{E}_\rho[L(h)]}.$$

Masegosa et al. have shown that the Chebyshev-Cantelli inequality is always at least as tight as second-order Markov's inequality (below we provide an alternative and more intuitive proof of this fact) and, therefore, the oracle bound in Theorem 2.4 is always at least as tight as the oracle bound in Theorem 2.2. However, the presence of $\mathbb{E}_{\rho^2}[L(h, h')]$ and $\mathbb{E}_\rho[L(h)]$ in the denominator make empirical estimation and optimization of the bound in Theorem 2.4 impractical, and Theorem 2.2 was the only practically applicable second-order bound so far.

2.4 Main Contributions

We present three main contributions: (1) a new form of the Chebyshev-Cantelli inequality, which is convenient for optimization; (2) an application of the inequality

to the analysis of weighted majority vote; and (3) a PAC-Bayes-Bennett inequality, which is used to bound the risk with an offset in the bound for weighted majority vote. We start with the new form of Chebyshev-Cantelli inequality, which can be seen as a refinement of second-order Markov's inequality or as an intermediate step in the proof of the Chebyshev-Cantelli inequality.

Theorem 2.5. *For any $\varepsilon > 0$ and all $\mu < \varepsilon$*

$$\mathbb{P}(Z \geq \varepsilon) \leq \frac{\mathbb{E}[(Z - \mu)^2]}{(\varepsilon - \mu)^2}.$$

Proof.

$$\mathbb{P}(Z \geq \varepsilon) = \mathbb{P}(Z - \mu \geq \varepsilon - \mu) \leq \mathbb{P}((Z - \mu)^2 \geq (\varepsilon - \mu)^2) \leq \frac{\mathbb{E}[(Z - \mu)^2]}{(\varepsilon - \mu)^2}. \quad \square$$

The inequality can also be written as

$$\mathbb{P}(Z \geq \varepsilon) \leq \frac{\mathbb{E}[(Z - \mu)^2]}{(\varepsilon - \mu)^2} = \frac{\mathbb{E}[Z^2] - 2\mu\mathbb{E}[Z] + \mu^2}{(\varepsilon - \mu)^2}. \quad (2.1)$$

The bound is minimized by $\mu^* = \mathbb{E}[Z] - \frac{\mathbb{V}[Z]}{\varepsilon - \mathbb{E}[Z]}$, which can be verified by taking a derivative of the bound with respect to μ . Note that μ^* can take negative values. Substitution of μ^* into the bound and simple algebraic manipulations recover the Chebyshev-Cantelli inequality, whereas $\mu = 0$ recovers second-order Markov's inequality. The main advantage of Theorem 2.5 over the Chebyshev-Cantelli inequality is ease of estimation and optimization due to absence of the variance term in the denominator.

Equation (2.1) leads to two new second-order oracle bounds for the weighted majority vote, given in Theorems 2.6 and 2.7.

Theorem 2.6. *In multiclass classification, for all ρ and all $\mu < 0.5$*

$$L(\text{MV}_\rho) \leq \frac{\mathbb{E}_{\rho^2}[L(h, h')] - 2\mu\mathbb{E}_\rho[L(h)] + \mu^2}{(0.5 - \mu)^2}.$$

Proof. As in the previous section, we take $Z = \mathbb{E}_\rho[\mathbf{1}(h(X) \neq Y)]$, so that $L(\text{MV}_\rho) \leq \mathbb{P}(Z \geq 0.5)$. The result follows by (2.1) and the calculations of $\mathbb{E}[Z^2]$ and $\mathbb{E}[Z]$ from the previous section. Note that the result is a deterministic statement. \square

For $\mu = 0$, Theorem 2.6 recovers Theorem 2.2, but if $\mu^* = \mathbb{E}_\rho[L(h)] - \frac{\mathbb{E}_{\rho^2}[L(h, h')] - \mathbb{E}[L(h)]^2}{0.5 - \mathbb{E}_\rho[L(h)]} \neq 0$, then substitution of μ^* into the theorem yields a tighter oracle bound. At the same time, substitution of μ^* recovers Theorem 2.4, but the great advantage of Theorem 2.6 is that the bound allows easy empirical estimation and optimization with respect to ρ , due to absence of $\mathbb{E}_{\rho^2}[L(h, h')]$ and $\mathbb{E}_\rho[L(h)]$ in the denominator. Thus, Theorem 2.6 has the oracle tightness of Theorem 2.4 and the ease of estimation and optimization of Theorem 2.2. The oracle quantities $\mathbb{E}_{\rho^2}[L(h, h')]$ and $\mathbb{E}_\rho[L(h)]$ can be bounded using PAC-Bayes-kl or PAC-Bayes- λ inequalities, as discussed in the next section.

In order to present the second oracle bound we introduce a new quantity. For a pair of hypotheses h and h' and a constant μ , we define *tandem loss with μ -offset*, for brevity *μ -tandem loss*, as

$$\ell_\mu(h(X), h'(X), Y) = (\mathbf{1}(h(X) \neq Y) - \mu)(\mathbf{1}(h'(X) \neq Y) - \mu). \quad (2.2)$$

Note that it can take negative values. We denote its expectation by

$$L_\mu(h, h') = \mathbb{E}_D[\ell_\mu(h(X), h'(X), Y)] = \mathbb{E}_D[(\mathbf{1}(h(X) \neq Y) - \mu)(\mathbf{1}(h'(X) \neq Y) - \mu)].$$

With $Z = \mathbb{E}_\rho[\mathbf{1}(h(X) \neq Y)]$ as before, we have

$$\begin{aligned} \mathbb{E}[(Z - \mu)^2] &= \mathbb{E}_D[(\mathbb{E}_\rho[\mathbf{1}(h(X) \neq Y) - \mu])^2] \\ &= \mathbb{E}_{\rho^2}[\mathbb{E}_D[(\mathbf{1}(h(X) \neq Y) - \mu)(\mathbf{1}(h'(X) \neq Y) - \mu)]] = \mathbb{E}_{\rho^2}[L_\mu(h, h')]. \end{aligned}$$

Now we present our second oracle bound.

Theorem 2.7. *In multiclass classification, for all ρ and all $\mu < 0.5$*

$$L(\text{MV}_\rho) \leq \frac{\mathbb{E}_{\rho^2}[L_\mu(h, h')]}{(0.5 - \mu)^2}.$$

Proof. The result follows by Theorem 2.5 and the calculation above. Note that the inequality is a deterministic statement. \square

In order to discuss the advantage of Theorem 2.7, we define the variance of the μ -tandem loss

$$\mathbb{V}_\mu(h, h') = \mathbb{E}_D[(\mathbf{1}(h(X) \neq Y) - \mu)(\mathbf{1}(h'(X) \neq Y) - \mu) - L_\mu(h, h')]^2.$$

If the variance of the μ -tandem loss is small, we can use Bernstein-type inequalities to obtain tighter estimates compared to kl-type inequalities.

We bound the μ -tandem loss using our next contribution, the PAC-Bayes-Bennett inequality, which improves on the PAC-Bayes-Bernstein inequality derived by Seldin et al. (2012) and may be of independent interest. The inequality holds for any loss function with bounded length of the range, we use $\tilde{\ell}$ and matching tilde-marked quantities to distinguish it from the zero-one loss ℓ . We let $\tilde{L}(h) = \mathbb{E}_D[\tilde{\ell}(h(X), Y)]$ and $\tilde{V}(h) = \mathbb{E}_D[(\tilde{\ell}(h(X), Y) - \tilde{L}(h))^2]$ be the expected tilde-loss of h and its variance and let $\hat{\tilde{L}}(h, S) = \frac{1}{n} \sum_{i=1}^n \tilde{\ell}(h(X_i), Y_i)$ be the empirical tilde-loss of h on a sample S .

Theorem 2.8 (PAC-Bayes-Bennett inequality). *Let $\tilde{\ell}(\cdot, \cdot)$ be an arbitrary loss function taking values in an interval of length b , and assume that $\tilde{V}(h)$ is finite for all h . Let $\phi(x) = e^x - x - 1$. Then for any distribution π on \mathcal{H} that is independent of S and any $\gamma > 0$ and $\delta \in (0, 1)$, with probability at least $1 - \delta$ over a random draw of S , for all distributions ρ on \mathcal{H} simultaneously:*

$$\mathbb{E}_\rho[\tilde{L}(h)] \leq \mathbb{E}_\rho[\hat{\tilde{L}}(h, S)] + \frac{\phi(\gamma b)}{\gamma b^2} \mathbb{E}_\rho[\tilde{V}(h)] + \frac{\text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta}}{\gamma n}.$$

The proof is based on a change of measure argument combined with Bennett's inequality, the details are provided in Appendix 2.8.1. Note that the result holds for a fixed (but arbitrary) $\gamma > 0$. In case of optimization with respect to γ a union bound has to be applied. For a fixed ρ the bound is convex in γ and for a fixed γ it is convex in ρ , although it is not necessarily jointly convex in ρ and γ . See Appendix 2.8.4 for optimization details. The PAC-Bayes-Bennett inequality is identical to the PAC-Bayes-Bernstein inequality of Seldin et al. (2012, Theorem 7), except that in the latter the coefficient in front of $\mathbb{E}_\rho[\tilde{V}[h]]$ is $(e - 2)\gamma$ instead of $\frac{\phi(\gamma b)}{\gamma b^2}$. The result improves on the result of Seldin et al. in two ways. First, in the result of Seldin et al. γ is restricted to the $(0, 1/b]$ interval, whereas in our result γ is unrestricted from above. And second, we can rewrite the coefficient in front of the variance as $\frac{\phi(\gamma b)}{\gamma b^2} = \frac{\phi(\gamma b)}{\gamma^2 b^2} \gamma$, where $\frac{\phi(\gamma b)}{\gamma^2 b^2}$ is a monotonically increasing function of γ , which in the interval $\gamma \in (0, 1/b]$ satisfies $\lim_{\gamma \rightarrow 0} \frac{\phi(\gamma b)}{\gamma^2 b^2} = \frac{1}{2}$ and for $\gamma = 1/b$ it gives $\frac{\phi(\gamma b)}{\gamma^2 b^2} = (e - 2)$. Thus, PAC-Bayes-Bennett is always at least as tight as PAC-Bayes-Bernstein and, at the same time, for $\gamma < 1/b$ it improves the constant coefficient in front of the variance from $(e - 2) \approx 0.72$ down to 0.5 for $\gamma \rightarrow 0$. For $\gamma > 1/b$ PAC-Bayes-Bennett also improves on PAC-Bayes-Bernstein, because PAC-Bayes-Bernstein uses the suboptimal value $\gamma = 1/b$ dictated by its restricted range of γ .

2.5 From oracle to empirical bounds

We obtain empirical bounds on the oracle quantities $\mathbb{E}_{\rho^2}[L(h, h')]$ and $\mathbb{E}_\rho[L(h)]$ in Theorem 2.6 and $\mathbb{E}_{\rho^2}[L_\mu(h, h')]$ in Theorem 2.7 by using PAC-Bayesian inequalities. The empirical counterpart of the expected tandem loss is the empirical tandem loss

$$\hat{L}(h, h', S) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(h(X_i) \neq Y_i) \mathbf{1}(h'(X_i) \neq Y_i).$$

For bounding $\mathbb{E}_{\rho^2}[L(h, h')]$ and $\mathbb{E}_\rho[L(h)]$ we use either PAC-Bayes-kl or PAC-Bayes- λ inequalities, both cited below. We use $\text{KL}(\rho||\pi)$ to denote the Kullback-Leibler divergence between distributions ρ and π on \mathcal{H} and $\text{kl}(p||q)$ to denote the Kullback-Leibler divergence between two Bernoulli distributions with biases p and q .

Theorem 2.9 (PAC-Bayes-kl Inequality, Seeger, 2002, Maurer, 2004). *For any probability distribution π on \mathcal{H} that is independent of S and any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over a random draw of a sample S , for all distributions ρ on \mathcal{H} simultaneously:*

$$\text{kl} \left(\mathbb{E}_\rho[\hat{L}(h, S)] \middle\| \mathbb{E}_\rho[L(h)] \right) \leq \frac{\text{KL}(\rho||\pi) + \ln(2\sqrt{n}/\delta)}{n}. \quad (2.3)$$

Theorem 2.10 (PAC-Bayes- λ Inequality, Thiemann et al., 2017; Masegosa et al., 2020). *For any probability distribution π on \mathcal{H} that is independent of S and any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over a random draw of a sample S , for all distributions ρ on \mathcal{H} and all $\lambda \in (0, 2)$ and $\gamma > 0$ simultaneously:*

$$\mathbb{E}_\rho[L(h)] \leq \frac{\mathbb{E}_\rho[\hat{L}(h, S)]}{1 - \frac{\lambda}{2}} + \frac{\text{KL}(\rho||\pi) + \ln(2\sqrt{n}/\delta)}{\lambda \left(1 - \frac{\lambda}{2}\right) n}, \quad (2.4)$$

$$\mathbb{E}_\rho[L(h)] \geq \left(1 - \frac{\gamma}{2}\right) \mathbb{E}_\rho[\hat{L}(h, S)] - \frac{\text{KL}(\rho||\pi) + \ln(2\sqrt{n}/\delta)}{\gamma n}. \quad (2.5)$$

(The upper bound (2.4) is due to Thiemann et al. (2017) and the lower bound (2.5) is due to Masegosa et al. (2020), and the two bounds hold simultaneously.) The PAC-Bayes- λ inequality is an optimization-friendly relaxation of the PAC-Bayes-kl inequality. Therefore, for optimization of ρ we use the PAC-Bayes- λ inequality, the upper bound for $\mathbb{E}_{\rho^2}[L(h, h')]$ and the lower or upper bound for $\mathbb{E}_\rho[L(h)]$, depending

on the positiveness of μ , but once we have converged to a solution we use PAC-Bayes-kl to compute the final bound. The kl form provides both an upper and a lower bound through the upper and lower inverse of the kl.¹ Taking the oracle bound from Theorem 2.6 and bounding the oracle quantities using Theorem 3.10 we obtain the following result.

Theorem 2.11. *For any distribution π on \mathcal{H} that is independent of S , and any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over a random draw of S , for all distributions ρ on \mathcal{H} , and all μ, λ , and γ in the ranges specified below simultaneously, we have:*

- For $\mu \in [0, 0.5)$, $\lambda \in (0, 2)$, and $\gamma > 0$:

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \mu)^2} \left[\frac{\mathbb{E}_{\rho^2}[\hat{L}(h, h', S)]}{1 - \frac{\lambda}{2}} + \frac{2 \text{KL}(\rho \|\pi) + \ln(4\sqrt{n}/\delta)}{\lambda (1 - \frac{\lambda}{2}) n} - 2\mu \left(\left(1 - \frac{\gamma}{2}\right) \mathbb{E}_\rho[\hat{L}(h, S)] - \frac{\text{KL}(\rho \|\pi) + \ln(4\sqrt{n}/\delta)}{\gamma n} \right) + \mu^2 \right].$$

- For $\mu < 0$, $\lambda \in (0, 2)$, and $\gamma \in (0, 2)$:

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \mu)^2} \left[\frac{\mathbb{E}_{\rho^2}[\hat{L}(h, h', S)]}{1 - \frac{\lambda}{2}} + \frac{2 \text{KL}(\rho \|\pi) + \ln(4\sqrt{n}/\delta)}{\lambda (1 - \frac{\lambda}{2}) n} - 2\mu \left(\frac{\mathbb{E}_\rho[\hat{L}(h, S)]}{1 - \frac{\gamma}{2}} + \frac{\text{KL}(\rho \|\pi) + \ln(4\sqrt{n}/\delta)}{\gamma (1 - \frac{\gamma}{2}) n} \right) + \mu^2 \right].$$

Proof. The result follows by substitution of the upper bound (2.4) on $\mathbb{E}_{\rho^2}[L(h, h')]$ and the lower bound (2.5) on $\mathbb{E}_\rho[L(h)]$ in the case of positive μ , or the upper bound (2.4) on $\mathbb{E}_\rho[L(h)]$ in the case of negative μ , into Theorem 2.6. We note that $\text{KL}(\rho^2 \|\pi^2) = 2 \text{KL}(\rho \|\pi)$ (Germain et al., 2015, Page 814), which gives the factor 2 in front of the first KL term. The factor 4 in the logarithms comes from a union bound over the bounds on $\mathbb{E}_{\rho^2}[L(h, h')]$ and $\mathbb{E}_\rho[L(h)]$. \square

We note that both the loss and the tandem loss are Bernoulli random variables, and for Bernoulli random variables the PAC-Bayes-kl inequality is tighter than the PAC-Bayes-Bennett (Tolstikhin and Seldin, 2013). However, the empirical counterpart of

¹Reeb et al. (2018) and Letarte et al. (2019) provide alternative ways of direct minimization of the upper bound on $\mathbb{E}_\rho[L(h)]$ given by the upper inverse of kl in the PAC-Bayes-kl bound. We use the PAC-Bayes- λ relaxation due to its simplicity, and because it provides an easy way of simultaneous optimization of an upper bound on $\mathbb{E}_{\rho^2}[L(h, h')]$ and a lower or upper bound on $\mathbb{E}_\rho[L(h)]$ (depending on μ).

the expected μ -tandem loss is the empirical μ -tandem loss

$$\hat{L}_\mu(h, h', S) = \frac{1}{n} \sum_{i=1}^n (\mathbf{1}(h(X_i) \neq Y_i) - \mu)(\mathbf{1}(h'(X_i) \neq Y_i) - \mu),$$

and the μ -tandem losses are not Bernoulli. Therefore, we use the PAC-Bayes-Bennett inequality, which provides an advantage if the variance of the μ -tandem losses happens to be small. The expected and empirical variance of the μ -tandem losses of a pair of hypotheses h and h' are, respectively, defined by

$$\begin{aligned} \mathbb{V}_\mu(h, h') &= \mathbb{E}_D[(\mathbf{1}(h(X) \neq Y) - \mu)(\mathbf{1}(h'(X) \neq Y) - \mu) - L_\mu(h, h')]^2, \\ \hat{\mathbb{V}}_\mu(h, h', S) &= \frac{1}{n-1} \sum_{i=1}^n \left((\mathbf{1}(h(X_i) \neq Y_i) - \mu)(\mathbf{1}(h'(X_i) \neq Y_i) - \mu) - \hat{L}_\mu(h, h', S) \right)^2. \end{aligned}$$

The empirical variance $\hat{\mathbb{V}}_\mu(h, h', S)$ is an unbiased estimate of $\mathbb{V}_\mu(h, h')$.

Since the PAC-Bayes-Bennett inequality is stated in terms of the oracle variance $\mathbb{E}_\rho[\tilde{\mathbb{V}}(h)]$, we use the result by Tolstikhin and Seldin (2013, Equation (15)) to bound it in terms of the empirical variance. For a general loss function $\tilde{\ell}(\cdot, \cdot)$ (not necessarily within $[0, 1]$), we define the empirical variance of the loss of h by $\hat{\tilde{\mathbb{V}}}(h, S) = \frac{1}{n-1} \sum_{i=1}^n (\tilde{\ell}(h(X_i), Y_i) - \tilde{L}(h))^2$. We recall that \tilde{L} , $\tilde{\mathbb{V}}$, and $\hat{\tilde{L}}$ were defined above Theorem 2.8. We note that the result of Tolstikhin and Seldin assumes that the losses are bounded in the $[0, 1]$ interval. Rescaling to a general range introduces the squared range factor c^2 in front of the last term in the inequality below, since scaling a random variable by c scales the variance by c^2 .

Theorem 2.12 (Tolstikhin and Seldin, 2013). *Let $\tilde{\ell}(\cdot, \cdot)$ be an arbitrary bounded loss function and let c be the length of the loss range. Then for any distribution π on \mathcal{H} that is independent of S , any $\lambda \in \left(0, \frac{2(n-1)}{n}\right)$, and any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over a random draw of the sample S , for all distributions ρ on \mathcal{H} simultaneously:*

$$\mathbb{E}_\rho[\tilde{\mathbb{V}}(h)] \leq \frac{\mathbb{E}_\rho[\hat{\tilde{\mathbb{V}}}(h, S)]}{1 - \frac{\lambda n}{2(n-1)}} + \frac{c^2 (\text{KL}(\rho||\pi) + \ln \frac{1}{\delta})}{n\lambda \left(1 - \frac{\lambda n}{2(n-1)}\right)}.$$

We note that, similar to the PAC-Bayes-Bennett inequality, but in contrast to the PAC-Bayes- λ inequality, the inequality above holds for a fixed value of λ and in case of optimization over λ a union bound has to be applied.

The last thing that is left is to bound the length of the range of μ -tandem losses defined in equation (2.2).

Lemma 2.2. *For $\mu < 0.5$ we have that the length of the range of $\ell_\mu(\cdot, \cdot, \cdot)$ is $K_\mu = \max\{1 - \mu, 1 - 2\mu\}$.*

A proof is provided in Appendix 2.8.2. Taking together the results of Theorems 2.7, 2.8, 2.12, and Lemma 2.2 we obtain the following result.

Theorem 2.13. *For any parameter grid $\{\gamma_1, \dots, \gamma_{k_\gamma}\}$ and $\{\lambda_1, \dots, \lambda_{k_\lambda}\}$, where $\gamma_i > 0$ for all i and $\lambda_i \in \left(0, \frac{2(n-1)}{n}\right)$ for all i , any distribution π on \mathcal{H} that is independent of S , and any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over a random draw of S , for all values of $\mu < 0.5$, all distributions ρ on \mathcal{H} , and all values of γ and λ in the parameter grid simultaneously:*

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \mu)^2} \left(\mathbb{E}_{\rho^2}[\hat{L}_\mu(h, h', S)] + \frac{2 \text{KL}(\rho \|\pi) + \ln \frac{2k_\gamma k_\lambda}{\delta}}{\gamma n} \right. \\ \left. + \frac{\phi(\gamma K_\mu)}{\gamma K_\mu^2} \left(\frac{\mathbb{E}_{\rho^2}[\hat{V}_\mu(h, h', S)]}{1 - \frac{\lambda n}{2(n-1)}} + \frac{K_\mu^2 \left(2 \text{KL}(\rho \|\pi) + \ln \frac{2k_\gamma k_\lambda}{\delta} \right)}{n\lambda \left(1 - \frac{\lambda n}{2(n-1)} \right)} \right) \right).$$

Proof. The result follows by reverse substitution of the result of Lemma 2.2 into Theorem 2.12, then into Theorem 2.8, and finally into Theorem 2.7. Since $\text{KL}(\rho^2 \|\pi^2) = 2 \text{KL}(\rho \|\pi)$, we have factor 2 in front of the KL terms. The factor $2k_\gamma k_\lambda$ comes from a union bound over the parameter grid and the bounds in Theorems 2.8 and 2.12. \square

2.6 Experiments

We start with a simulated comparison of the oracle bounds and then present an empirical evaluation on real data. The python source code for replicating the experiments is available at Github².

Comparison of the oracle bounds

Figure 2.1 depicts a comparison of the second-order oracle bound based on the Chebyshev-Cantelli inequality (Theorems 2.4, 2.6 and 2.7, which, as oracle bounds, are equivalent) and the second-order oracle bound based on the second-order Markov's

²<https://github.com/StephanLorenzen/MajorityVoteBounds>

inequality (Theorem 2.2). We plot the ratio of the right hand side of the bound in

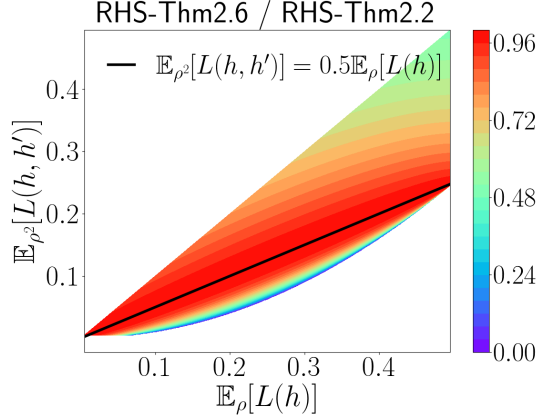


Figure 2.1: Theorem 2.6 vs. Theorem 2.2

Theorem 2.6 for the optimal value $\mu^* = \mathbb{E}_\rho[L(h)] - \frac{\mathbb{E}_{\rho^2}[L(h, h')] - \mathbb{E}_\rho[L(h)]^2}{0.5 - \mathbb{E}_\rho[L(h)]}$ to the value of the right hand side of the bound in Theorem 2.2. A simple calculation shows that if $\mathbb{E}_{\rho^2}[L(h, h')] = 0.5\mathbb{E}_\rho[L(h)]$, then $\mu^* = 0$, which recovers the bound in Theorem 2.2. The line $\mathbb{E}_{\rho^2}[L(h, h')] = 0.5\mathbb{E}_\rho[L(h)]$ is shown in black in Figure 2.1. We also note that $\mathbb{E}_\rho[L(h)]^2 \leq \mathbb{E}_{\rho^2}[L(h, h')] \leq \mathbb{E}_\rho[L(h)]$, which defines the feasible region in Figure 2.1. Whenever $\mathbb{E}_{\rho^2}[L(h, h')] \neq 0.5\mathbb{E}_\rho[L(h)]$ the Chebyshev-Cantelli inequality provides an improvement over second-order Markov's inequality. The region above the black line, where $\mathbb{E}_{\rho^2}[L(h, h')] > 0.5\mathbb{E}_\rho[L(h)]$, is the region of high correlation of errors and in this case majority vote yields little improvement over individual classifiers. In this region the first-order oracle bound is tighter than the second-order oracle bounds (see Appendix 2.8.3). The region below the black line, where $\mathbb{E}_{\rho^2}[L(h, h')] < 0.5\mathbb{E}_\rho[L(h)]$, is the region of low correlation of errors. In this region the second-order oracle bounds are tighter than the first-order oracle bound. Note that the potential for improvement below the black line is much higher than above it.

Empirical evaluation on real datasets

We studied the empirical performance of the bounds using standard random forest (Breiman, 2001) and a combination of heterogeneous classifiers on a subset of data sets from UCI and LibSVM repositories (Dua and Graff, 2019; Chang and Lin, 2011). An overview of the data sets is given in Appendix 2.8.5.1. The number of points varied from 3000 to 70000 with dimensions $d < 1000$. For each data set, we set aside 20% of the data for the test set S_{test} and used the remaining data S for ensemble

construction, weight optimization and bound evaluation. We evaluate the classifiers and bounds obtained by minimizing the tandem bound TND (Masegosa et al., 2020, Theorem 9), which is the empirical bound on the oracle tandem bound in Theorem 2.2, the Chebyshev-Cantelli bound with TND empirical loss estimate bound CCTND (Theorem 2.11), and the Chebyshev-Cantelli bound with PAC-Bayes-Bennett loss estimate bound CCPBB (Theorem 2.13). We made 10 repetitions of each experiment.

Ensemble construction and minimization of the bounds. We follow the construction used by Masegosa et al. (2020). The idea is to generate multiple random splits of the data set S into pairs of subsets $S = T_h \cup S_h$, such that $T_h \cap S_h = \emptyset$. Each hypothesis is trained on T_h and the empirical loss on S_h provides an unbiased estimate of its expected loss. Note that the splits cannot depend on the data. For our experiments, we generate these splits by bagging, where out-of-bag (OOB) samples S_h provide unbiased estimates of expected losses of individual hypotheses h . The resulting set of hypotheses produces an ensemble. As in the work of Masegosa et al., two modifications are required to apply the bounds: the empirical losses $\hat{L}(h, S)$ in the bounds are replaced by the validation losses $\hat{L}(h, S_h)$, and the sample size n is replaced by the minimal validation size $\min_h |S_h|$. For pairs of hypotheses (h, h') , we take the overlaps of their validation sets $S_h \cap S_{h'}$ to calculate an unbiased estimate of their tandem loss $\hat{L}(h, h', S_h \cap S_{h'})$, μ -tandem loss $\hat{L}_\mu(h, h', S_h \cap S_{h'})$, and the variance of the μ -tandem loss $\hat{V}_\mu(h, h', S_h \cap S_{h'})$, which replaces the corresponding empirical losses in the bounds. The sample size is then replaced by $\min_{h, h'} |S_h \cap S_{h'}|$. The details on bound minimization are provided in Appendix 2.8.4.

Optimizing weighted random forest. In the first experiment we compare TND, CCTND, and CCPBB bounds in the setting studied by Masegosa et al. (2020). We take 100 fully grown trees, use the Gini criterion for splitting, and consider \sqrt{d} features in each split. Figure 2.2a compares the loss of the random forest on S_{test} using either uniform weighting ρ_u or optimized weighting ρ^* found by minimization of the three bounds (we exclude the first-order bound from the comparison, since it was shown by Masegosa et al. that it significantly deteriorates the test error of the ensemble). While CCTND often performs similar to TND, we find that optimizing using CCPBB often improves accuracy. Figure 2.2b compares the tightness of the optimized CCTND and CCPBB bounds to the optimized TND bound. The CCTND is generally comparable to TND, while CCPBB is consistently looser than TND, mainly due to the union bounds. The numerical values for the losses and the bounds can be found in Tables 2.2 and 2.3 in Appendix 2.8.5.2.

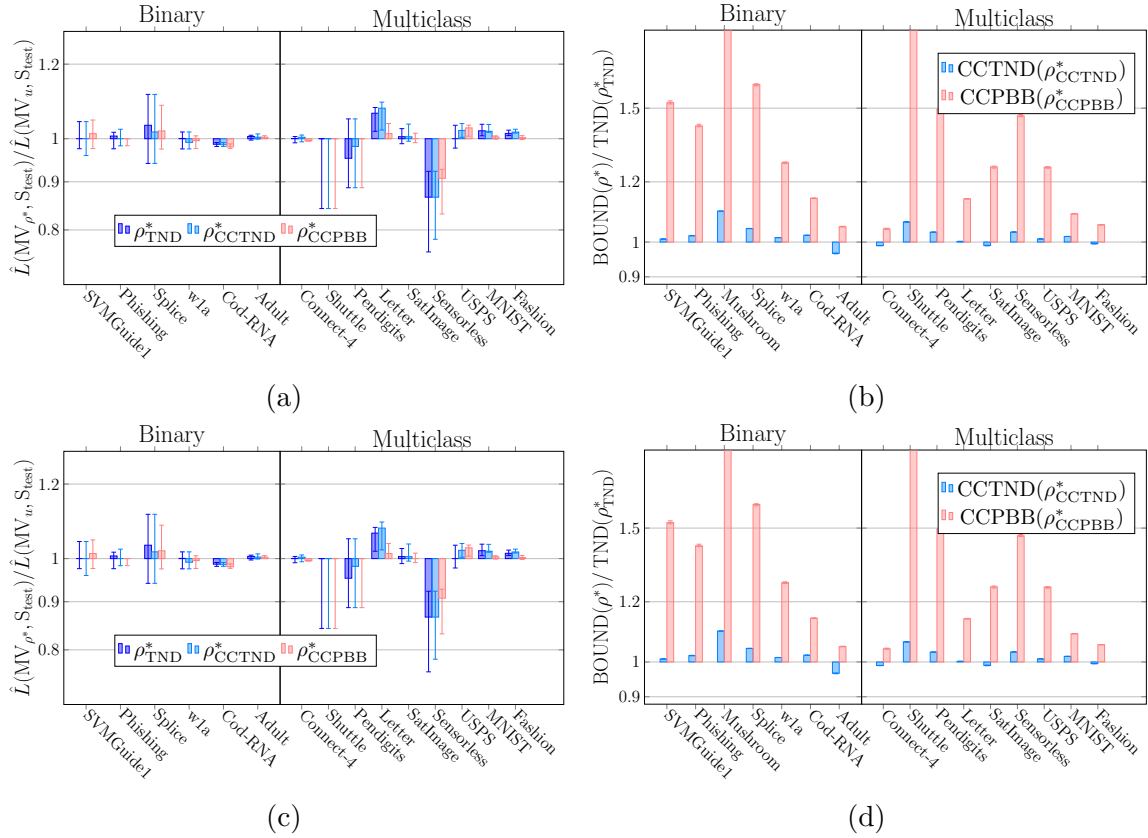


Figure 2.2: **(a,b) Optimized weighted random forest. (c,d) Ensembles with heterogeneous classifiers.** The median, 25%, and 75% quantiles of: (a,c) the ratio $\hat{L}(MV_{\rho^*}, S_{\text{test}}) / \hat{L}(MV_u, S_{\text{test}})$ of the test loss of the majority vote with optimized weighting ρ^* generated by TND, CCTND and CCPBB to the test loss of majority vote with uniform weighting, and (b,d) the ratio $\text{BOUND}(\rho^*) / \text{TND}(\rho_{\text{TND}}^*)$ of the CCTND and CCPBB bounds to the TND bound with the corresponding optimized weighting. The plots are on a logarithmic scale. Values above 1 represent degradation and values below 1 represent improvement. Data sets with $L(MV_u, S_{\text{test}}) = 0$ are left out in (a,c).

Ensembles with heterogeneous classifiers. In the second experiment, we consider ensembles of heterogeneous classifiers (Linear Discriminant Analysis, k -Nearest Neighbors, Decision Tree, Logistic Regression, and Gaussian Naive Bayes). A detailed description is provided in Appendix 2.8.5.3. Compared to random forests, the variation in performance of ensemble members is larger here. Figure 2.2c compares

the ratio of the loss of the majority vote with optimized weighting to the loss of majority vote with uniform weighting on S_{test} for ρ^* found by minimization of the first-order bound (FO), TND, CCTND, and CCPBB. The numerical values are given in Table 2.5 in Appendix 2.8.5.3. We observed that optimizing the FO tends to improve the ensemble accuracy in some cases but degrade in others. However, TND, CCTND, and CCPBB almost always improve the performance w.r.t. the uniform weighting. Table 2.5 also shows that choosing the best single hypothesis gives almost identical results as optimizing FO. Figure 2.2d compares the tightness of the CCTND and CCPBB bounds relative to the TND bound. The numerical values are given in Table 2.6 in Appendix 2.8.5.3. In this case, we have that CCTND is usually tighter than TND, while CCPBB is usually looser than TND due to the union bounds.

2.7 Discussion

We derived an optimization-friendly form of the Chebyshev-Cantelli inequality and applied it to derive two new forms of second-order oracle bounds for the weighted majority vote. The new oracle bounds bridge between the C-bounds (Germain et al., 2015) and the tandem bound (Masegosa et al., 2020) and take the best of both: the tightness of the Chebyshev-Cantelli inequality and the optimization and estimation convenience of the tandem bound. We also derived the PAC-Bayes-Bennett inequality, improving on the PAC-Bayes-Bernstein inequality of Seldin et al. (2012).

Our paper opens several directions for future research. One of them is a better treatment of parameter search in parametric bounds that would give tighter bounds than a union bound over a grid. It would also be interesting to find other applications for the new form of Chebyshev-Cantelli inequality and the PAC-Bayes-Bennett inequality.

2.8 Appendix

2.8.1 A proof of the PAC-Bayes-Bennett inequality (Theorem 2.8) and a comparison with the PAC-Bayes-Bernstein inequality

In this section we provide a proof of Theorem 2.8 and a numerical comparison with the PAC-Bayes-Bernstein inequality. The proof is based on the standard change of measure argument. We use the following version by Tolstikhin and Seldin (2013).

Lemma 2.3 (PAC-Bayes Lemma). *For any function $f_n : \mathcal{H} \times (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathbb{R}$ and for any distribution π on \mathcal{H} , such that π is independent of S , with probability at least $1 - \delta$ over a random draw of S , for all distributions ρ on \mathcal{H} simultaneously:*

$$\mathbb{E}_\rho[f_n(h, S)] \leq \text{KL}(\rho||\pi) + \ln \frac{1}{\delta} + \ln \mathbb{E}_\pi[\mathbb{E}_{S'}[e^{f_n(h, S')}]].$$

The second ingredient is Bennett's lemma, which is a bound on the moment generating function used in the proof of Bennett's inequality. Since we are unaware of a reference, we provide a proof below, which is essentially an intermediate step in the proof of Bennett's inequality (Boucheron et al., 2013, Theorem 2.9).

Lemma 2.4 (Bennett's Lemma). *Let $b > 0$ and let Z_1, \dots, Z_n be i.i.d. zero-mean random variables with finite variance, such that $Z_i \leq b$ for all i . Let $M_n = \sum_{i=1}^n Z_i$ and $V_n = \sum_{i=1}^n \mathbb{E}[Z_i^2]$. Let $\phi(u) = e^u - u - 1$. Then for any $\lambda > 0$:*

$$\mathbb{E} \left[e^{\lambda M_n - \frac{\phi(b\lambda)}{b^2} V_n} \right] \leq 1.$$

Proof. Since $u^{-2}\phi(u)$ is a non-decreasing function of $u \in \mathbb{R}$ (where at zero we continuously extend the function), for all $i \in [n]$ and $\lambda > 0$ we have

$$e^{\lambda Z_i} - \lambda Z_i - 1 \leq Z_i^2 \frac{\phi(b\lambda)}{b^2},$$

which implies

$$\mathbb{E} [e^{\lambda Z_i}] \leq 1 + \lambda \mathbb{E} [Z_i] + \frac{\phi(b\lambda)}{b^2} \mathbb{E} [Z_i^2] \leq e^{\frac{\phi(b\lambda)}{b^2} \mathbb{E}[Z_i^2]},$$

where the second inequality uses the assumption that $\mathbb{E}[Z_i] = 0$ and the fact that $1 + x \leq e^x$ for all $x \in \mathbb{R}$. By the above inequality and independence of the random variables,

$$\mathbb{E} \left[e^{\lambda M_n - \frac{\phi(b\lambda)}{b^2} V_n} \right] = \mathbb{E} \left[\prod_{i=1}^n e^{\lambda Z_i - \frac{\phi(b\lambda)}{b^2} \mathbb{E}[Z_i^2]} \right] = \prod_{i=1}^n \mathbb{E} \left[e^{\lambda Z_i - \frac{\phi(b\lambda)}{b^2} \mathbb{E}[Z_i^2]} \right] \leq 1.$$

□

Now we are ready to prove the theorem.

Proof of Theorem 2.8. We take $f_n(h, S) = \gamma n \left(\tilde{L}(h) - \hat{\tilde{L}}(h, S) \right) - \frac{\phi(\gamma b)}{b^2} n \tilde{V}(h)$. Then by Lemma 2.4 we have $\mathbb{E}_S[e^{f_n(h, S)}] \leq 1$. By plugging this into Lemma 2.3, normalizing by γn , and changing sides, we obtain the result. □

Numerical comparison of PAC-Bayes-Bennett and PAC-Bayes-Bernstein bound

Figure 2.3 provides a numerical comparison of PAC-Bayes-Bennett and PAC-Bayes-Bernstein inequalities (Theorem 2.8 and Theorem 7 by Tolstikhin and Seldin (2013)).

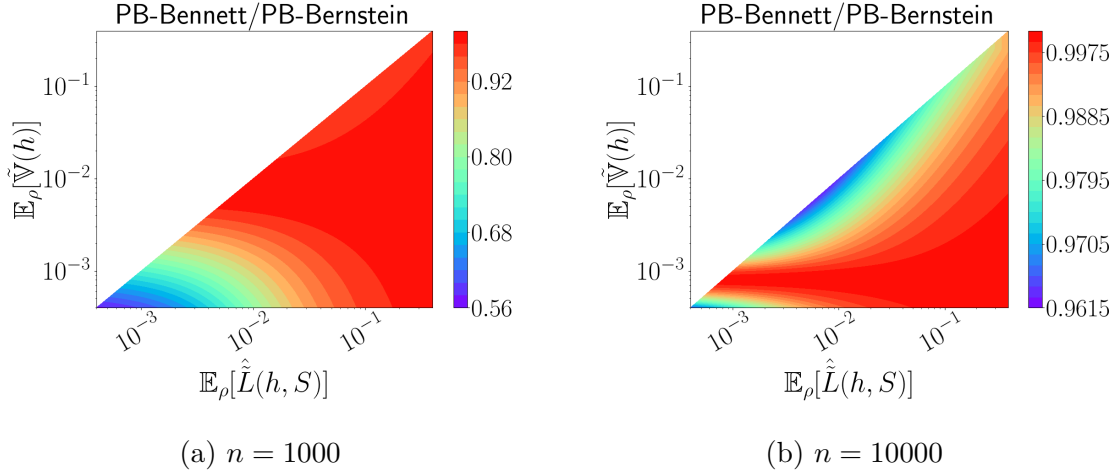


Figure 2.3: The ratio of PAC-Bayes Bennett to PAC-Bayes Bernstein bound as a function of $\mathbb{E}_\rho[\hat{L}(h, S)]$ and $\mathbb{E}_\rho[\tilde{V}(h)]$. We set $\text{KL}(\rho||\pi) = 5$ and $\delta = 0.05$. The value of n is provided in the captions of the subfigures.

2.8.2 Proof of Lemma 2.2

Proof. Recall that

$$\ell_\mu(h(X), h'(X), Y) = (\mathbb{1}(h(X) \neq Y) - \mu)(\mathbb{1}(h'(X) \neq Y) - \mu) \in \{(1 - \mu)^2, -\mu(1 - \mu), \mu^2\}.$$

For $\mu < 0.5$, we have $-\mu(1 - \mu) < (1 - \mu)^2$ and $\mu^2 < (1 - \mu)^2$. Therefore, $\ell_\mu(h(X), h'(X), Y) \leq (1 - \mu)^2$. Furthermore, for $\mu < 0$ we have $\mu^2 < -\mu(1 - \mu)$, and for $\mu > 0$ we have $-\mu(1 - \mu) \leq \mu^2$. Therefore, for $\mu < 0.5$ we have $\ell_\mu(h(X), h'(X), Y) \geq \min\{-\mu(1 - \mu), \mu^2\}$.

By combining the upper and the lower bound, we obtain

$$\begin{aligned} K_\mu &= (1 - \mu)^2 - \min\{-\mu(1 - \mu), \mu^2\} \\ &= \max\{(1 - \mu)^2 - (-\mu(1 - \mu)), (1 - \mu)^2 - \mu^2\} \\ &= \max\{1 - \mu, 1 - 2\mu\}. \end{aligned}$$

□

2.8.3 Comparison of the first and second-order oracle bounds

In this section we show that if $\mathbb{E}_\rho[L(h)] < 0.5$ and $\mathbb{E}_{\rho^2}[L(h, h')] > 0.5\mathbb{E}_\rho[L(h)]$, then the first-order oracle bound is tighter than the second-order oracle bounds, and if $\mathbb{E}_\rho[L(h)] < 0.5$ and $\mathbb{E}_{\rho^2}[L(h, h')] < 0.5\mathbb{E}_\rho[L(h)]$, then it is the other way around.

For comparison of the first-order oracle bound $L(\text{MV}_\rho) \leq 2\mathbb{E}_\rho[L(h)]$ vs. the second-order oracle tandem bound $L(\text{MV}_\rho) \leq 4\mathbb{E}_{\rho^2}[L(h, h')]$ the statement above is evident.

For the second-order oracle bounds based on the Chebyshev-Cantelli inequality we have

$$\begin{aligned} \frac{\mathbb{E}_{\rho^2}[L(h, h')] - \mathbb{E}_\rho[L(h)]^2}{0.25 + \mathbb{E}_{\rho^2}[L(h, h')] - \mathbb{E}_\rho[L(h)]} & \quad \text{vs.} \quad 2\mathbb{E}_\rho[L(h)], \\ \frac{\mathbb{E}_{\rho^2}[L(h, h')] - \mathbb{E}_\rho[L(h)]^2}{\mathbb{E}_{\rho^2}[L(h, h')](1 - 2\mathbb{E}_\rho[L(h)])} & \quad \text{vs.} \quad 0.5\mathbb{E}_\rho[L(h)] + 2\mathbb{E}_\rho[L(h)]\mathbb{E}_{\rho^2}[L(h, h')] - 2\mathbb{E}_\rho[L(h)]^2, \\ \frac{\mathbb{E}_{\rho^2}[L(h, h')]}{\mathbb{E}_{\rho^2}[L(h, h')]} & \quad \text{vs.} \quad 0.5\mathbb{E}_\rho[L(h)], \end{aligned}$$

where under the assumption that $\mathbb{E}_\rho[L(h)] < 0.5$ we can cancel $(1 - 2\mathbb{E}_\rho[L(h)])$, since it is positive, and the result is again evident.

2.8.4 Minimization of the bounds

In this section we provide technical details on minimization of the bounds in Theorems 2.11 and 2.13. As most of the other PAC-Bayesian works, we take π to be a union distribution over the hypotheses in both cases. As discussed in Section 2.6, we build a set of data-dependent hypotheses by splitting the data set S into pairs of subsets $S = T_h \cup S_h$, such that $T_h \cap S_h = \emptyset$, training h on T_h and calculating an unbiased loss estimate $\hat{L}(h, S_h)$ on S_h . For tandem losses we compute the unbiased estimates $\hat{L}(h, h', S_h \cap S_{h'})$ on the intersections of the corresponding sets S_h and $S_{h'}$.

2.8.4.1 Minimization of the bound in Theorem 2.11

The adjustment of the bound from Theorem 2.11 to this construction is for $\mu \geq 0$:

$$\begin{aligned} L(\text{MV}_\rho) \leq & \frac{1}{(0.5 - \mu)^2} \left[\frac{\mathbb{E}_{\rho^2}[\hat{L}(h, h', S_h \cap S_{h'})]}{1 - \frac{\lambda}{2}} + \frac{2 \text{KL}(\rho \parallel \pi) + \ln(4\sqrt{m}/\delta)}{\lambda \left(1 - \frac{\lambda}{2}\right) m} \right. \\ & \left. - 2\mu \left(\left(1 - \frac{\gamma}{2}\right) \mathbb{E}_\rho[\hat{L}(h, S_h)] - \frac{\text{KL}(\rho \parallel \pi) + \ln(4\sqrt{n}/\delta)}{\gamma n} \right) + \mu^2 \right], \end{aligned}$$

and for $\mu < 0$:

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \mu)^2} \left[\frac{\mathbb{E}_{\rho^2}[\hat{L}(h, h', S_h \cap S_{h'})]}{1 - \frac{\lambda}{2}} + \frac{2 \text{KL}(\rho \|\pi) + \ln(4\sqrt{m}/\delta)}{\lambda(1 - \frac{\lambda}{2})m} - 2\mu \left(\frac{\mathbb{E}_\rho[\hat{L}(h, S_h)]}{1 - \frac{\gamma}{2}} + \frac{\text{KL}(\rho \|\pi) + \ln(4\sqrt{n}/\delta)}{\gamma(1 - \frac{\gamma}{2})n} \right) + \mu^2 \right],$$

where $m = \min_{h, h'} |S_h \cap S_{h'}|$ and $n = \min_h |S_h|$. Below we provide the pseudocode and derive update rules for μ , λ , γ , and ρ for alternating minimization of this bound.

Algorithm 1: Minimization of the bound in Theorem 2.11

Input: m, n , tandem losses $\hat{L}(h, h', S_h \cap S_{h'})$ for all h, h' , and Gibbs losses $\hat{L}(h, S_h)$ for all h

Initialize: $\rho = \pi$ and $\mu = 0$

while The improvement of the bound is larger than 10^{-9} **do**

 Compute λ_ρ^* , the optimal λ given ρ

 Compute γ_ρ^* , the optimal γ given ρ and μ

 Compute the bound using ρ , μ , λ_ρ^* and γ_ρ^*

 Compute new μ_ρ^* , the optimal μ given ρ , λ_ρ^* and γ_ρ^*

 Update the new distribution ρ' with gradient descent given μ , λ_ρ^* and γ_ρ^*

 Let $\rho = \rho'$ and $\mu = \mu_\rho^*$

end while

Optimal λ given ρ Minimization of the bound with respect to λ is identical to minimization of the tandem bound by Masegosa et al. (2020, Theorem 9). Masegosa et al. derive the optimal value of λ :

$$\lambda_\rho^* = \frac{2}{\sqrt{\frac{2m\mathbb{E}_{\rho^2}[\hat{L}(h, h', S_h \cap S_{h'})]}{2 \text{KL}(\rho \|\pi) + \ln \frac{4\sqrt{m}}{\delta}} + 1 + 1}}.$$

Optimal γ given ρ and μ Minimization of the bound with respect to γ in the case of $\mu \geq 0$ is analogous to minimization of the bound by Masegosa et al. (2020, Theorem 10) with respect to γ . Masegosa et al. derive the optimal value of γ :

$$\gamma_\rho^* = \sqrt{\frac{2 \text{KL}(\rho \|\pi) + \ln(16n/\delta^2)}{n\mathbb{E}_\rho[\hat{L}(h, S_h)]}}.$$

On the other hand, the optimal γ in the case of $\mu < 0$ is analogous to the optimal λ above:

$$\gamma_\rho^* = \frac{2}{\sqrt{\frac{2n\mathbb{E}_\rho[\hat{L}(h, S_h)]}{\text{KL}(\rho||\pi) + \ln \frac{4\sqrt{n}}{\delta}} + 1 + 1}}.$$

Optimal μ given ρ Given ρ , we can compute the optimal λ_ρ^* and γ_ρ^* by the above formulas. Let

$$U_T(\rho) := \frac{\mathbb{E}_{\rho^2}[\hat{L}(h, h', S_h \cap S_{h'})]}{1 - \frac{\lambda_\rho^*}{2}} + \frac{2 \text{KL}(\rho||\pi) + \ln(4\sqrt{m}/\delta)}{\lambda_\rho^* \left(1 - \frac{\lambda_\rho^*}{2}\right) m},$$

$$L_G(\rho) := \begin{cases} \left(1 - \frac{\gamma_\rho^*}{2}\right) \mathbb{E}_\rho[\hat{L}(h, S_h)] - \frac{\text{KL}(\rho||\pi) + \ln(4\sqrt{n}/\delta)}{\gamma_\rho^* n}, & \mu \geq 0 \\ \frac{\mathbb{E}_\rho[\hat{L}(h, S_h)]}{1 - \frac{\gamma_\rho^*}{2}} + \frac{\text{KL}(\rho||\pi) + \ln(4\sqrt{n}/\delta)}{\gamma_\rho^* \left(1 - \frac{\gamma_\rho^*}{2}\right) n}, & \mu < 0 \end{cases}$$

Then the optimal μ is

$$\mu_\rho^* = \frac{\frac{1}{2}L_G(\rho) - U_T(\rho)}{\frac{1}{2} - L_G(\rho)}.$$

Gradient w.r.t. ρ given λ , γ and μ Minimization of the bound w.r.t. ρ is equivalent to constrained optimization of $f(\rho) = a\mathbb{E}_{\rho^2}[\hat{L}(h, h', S_h \cap S_{h'})] - 2b\mathbb{E}_\rho[\hat{L}(h, S_h)] + 2c\text{KL}(\rho||\pi)$, where for $\mu \geq 0$, $a = 1/(1 - \lambda/2)$, $b = \mu(1 - \gamma/2)$ and $c = 1/(\lambda(1 - \lambda/2)m) + \mu/(\gamma n)$, and for $\mu < 0$, $a = 1/(1 - \lambda/2)$, $b = \mu/(1 - \gamma/2)$, and $c = 1/(\lambda(1 - \lambda/2)m) - \mu/(\gamma(1 - \gamma/2)n)$. The constraint is that ρ is a probability distribution. We optimize ρ by projected gradient descent, where we iteratively take steps in the direction of the negative gradient of f and project the result onto the probability simplex.

We use \hat{L} to denote the vector of empirical losses and \hat{L}_{tnd} to denote the matrix of tandem losses. Let ∇f denote the gradient of f w.r.t. ρ and $(\nabla f)_h$ the h -th coordinate of the gradient. We have:

$$(\nabla f)_h = 2 \left(a \sum_{h'} \rho(h') \hat{L}(h, h', S_h \cap S_{h'}) - b \hat{L}(h, S_h) + c \left(1 + \ln \frac{\rho(h)}{\pi(h)} \right) \right),$$

$$\nabla f = 2 \left(a \hat{L}_{\text{tnd}} \rho - b \hat{L} + c \left(1 + \ln \frac{\rho}{\pi} \right) \right).$$

Gradient descent optimization w.r.t. ρ To optimize the weighting ρ , we applied iRProp+ for the gradient based optimization, a first-order method with adaptive individual step sizes (Igel and Hüsken, 2003; Florescu and Igel, 2018), until the bound did not improve for 10 iterations.

2.8.4.2 Minimization of the bound in Theorem 2.13

We start with the details on construction of the grid of μ , λ and γ .

The μ grid for Theorem 2.13

We were unable to find a closed-form solution for minimization of the bound w.r.t. μ and applied a heuristic. Empirically we observed that the bound was quasiconvex in μ (we were unable to prove that it is always the case) and applied binary search for μ in the grid. Note that even if we take a grid of μ , we don't need a union bound since the bound holds with high probability for all μ simultaneously.

We then consider the relevant range of μ . By Theorem 2.5, we have $\mu < 0.5$. At the same time, $\mu^* = \frac{0.5\mathbb{E}_\rho[L(h)] - \mathbb{E}_{\rho^2}[L(h,h)]}{0.5 - \mathbb{E}_\rho[L(h)]}$, and in Section 2.6 we have shown that the primary region of interest is where $\mathbb{E}_{\rho^2}[L(h, h')] < 0.5\mathbb{E}_\rho[L(h)]$, which corresponds to $\mu^* > 0$. However, since $\mathbb{E}_{\rho^2}[L(h, h)]$ and $\mathbb{E}_\rho[L(h)]$ are unobserved and we use an upper bound for the first and a lower bound for the second instead, we take a broader range of μ . By making a mild assumption that the upper bound for the tandem loss $\mathbb{E}_{\rho^2}[L(h, h')]$ is at most 0.25 and the lower bound for the Gibbs loss $\mathbb{E}_\rho[L(h)]$ is at most 0.5, we have $\mu \in [-0.5, 0.5)$. We take 400 uniformly spaced points in the selected range for the CCPBB bound.

The λ grid for Theorem 2.13

The parameter λ comes from Theorem 2.12. The theorem is identical to the result by Tolstikhin and Seldin (2013, Equation (15)), except rescaling, but rescaling happens on top of the bound and has no effect on the λ -grid. Therefore, we use the grid proposed by Tolstikhin and Seldin. Namely, we take

$$\lambda_i = c_1^{i-1} \frac{2(n-1)}{n} \left(\sqrt{\frac{n-1}{\ln(1/\delta_1)} + 1} + 1 \right)^{-1}$$

for $i \in \{1, \dots, k_\lambda\}$ and

$$k_\lambda = \left\lceil \frac{1}{\ln c_1} \ln \left(\frac{1}{2} \sqrt{\frac{n-1}{\ln(1/\delta_1)} + 1} + \frac{1}{2} \right) \right\rceil.$$

In the experiments we took $c_1 = 1.05$ and $\delta_1 = \delta/2$.

The γ grid for Theorem 2.13

The parameter γ comes from Theorem 2.8. By taking the first two derivatives we can verify that for a fixed ρ the PAC-Bayes-Bennett bound is convex in γ and at the minimum point the optimal value of γ satisfies

$$e^{(\gamma_\rho^* b - 1)} (\gamma_\rho^* b - 1) = \frac{1}{e} \left(\frac{b^2 \left(\text{KL}(\rho \|\pi) + \ln \frac{1}{\delta_2} \right)}{n \mathbb{E}_\rho[\tilde{\mathbb{V}}(h)]} - 1 \right).$$

Thus, the optimal value of γ is given by

$$\gamma_\rho^* = \frac{1}{b} \left(W_0 \left(\frac{1}{e} \left(\frac{b^2 \left(\text{KL}(\rho \|\pi) + \ln \frac{1}{\delta_2} \right)}{n \mathbb{E}_\rho[\tilde{\mathbb{V}}(h)]} - 1 \right) \right) + 1 \right),$$

where W_0 is the principal branch of the Lambert W function, which is defined as the inverse of the function $f(x) = xe^x$.

In order to define a grid for γ we first determine the relevant range for γ_ρ^* . We note that the variance $\mathbb{E}_\rho[\tilde{\mathbb{V}}(h)]$ is estimated using Theorem 2.12, which assumes that the length of the range of the loss $\tilde{\ell}(\cdot, \cdot)$ is c . The loss range provides a trivial upper bound on the variance $\mathbb{E}_\rho[\tilde{\mathbb{V}}(h)] \leq \frac{c^2}{4}$. At the same time, we have $\lambda \left(1 - \frac{\lambda n}{2(n-1)} \right) \leq \frac{n-1}{2n}$ (it is a downward-pointing parabola) and, therefore, the right hand side of the bound in Theorem 2.12 is at least the value of its second term, which is at least $\frac{2c^2 \ln \frac{1}{\delta_1}}{n-1}$, since $\text{KL}(\rho \|\pi) \geq 0$. Thus, we obtain that the estimate of $\mathbb{E}_\rho[\tilde{\mathbb{V}}(h)]$ is in the range $\left[\frac{2c^2 \ln \frac{1}{\delta_1}}{n-1}, \frac{c^2}{4} \right]$. We use $V_{\min} = \frac{2c^2 \ln \frac{1}{\delta_1}}{n-1}$ to denote the lower bound of this range.

Since $W_0(\cdot)$ is a monotonically increasing function, $\text{KL}(\rho \|\pi) \geq 0$, and the estimate of

$\mathbb{E}_\rho[\tilde{\mathbb{V}}(h)]$ is at most $\frac{c^2}{4}$, we obtain that γ_ρ^* satisfies

$$\begin{aligned}\gamma_\rho^* &= \frac{1}{b} \left(W_0 \left(\frac{1}{e} \left(\frac{b^2 \left(\text{KL}(\rho||\pi) + \ln \frac{1}{\delta_2} \right)}{n\mathbb{E}_{\rho^2}[\tilde{\mathbb{V}}(h)]} - 1 \right) \right) + 1 \right) \\ &\geq \frac{1}{b} \left(W_0 \left(\frac{1}{e} \left(\frac{4b^2}{nc^2} \ln \frac{1}{\delta_2} - 1 \right) \right) + 1 \right) \stackrel{def}{=} \gamma_{\min}.\end{aligned}$$

For an upper bound we observe that since $\mathbb{E}_\rho[\tilde{L}(h)] - \mathbb{E}_\rho[\hat{L}(h, S)]$ is trivially bounded by b , the bound in Theorem 2.12 is only interesting if it is smaller than b and, in particular, $\frac{\phi(\gamma b)}{\gamma b^2} \mathbb{E}_\rho[\tilde{\mathbb{V}}(h)] \leq b$.

This gives

$$b \geq \frac{\phi(\gamma b)}{\gamma b^2} \mathbb{E}_\rho[\tilde{\mathbb{V}}(h)] \geq \frac{\phi(\gamma b)}{\gamma b^2} V_{\min}.$$

Thus, γ should satisfy

$$\phi(\gamma b) \leq \frac{\gamma b^3}{V_{\min}},$$

which gives that the maximal value of γ , denoted γ_{max} , is the positive root of

$$H(\gamma) = e^{\gamma b} - \gamma b \left(1 + \frac{b^2}{V_{\min}} \right) - 1 = 0.$$

Let $\alpha = (1 + b^2/V_{\min})^{-1} \in (0, 1)$, and $x = -\gamma b - \alpha$. Then the above problem is equivalent to finding the root of $f(x) = xe^x - d$ for $d = -\alpha e^{-\alpha}$, which can again be solved by applying the Lambert W function. Since for $\alpha \in (0, 1)$, we have $d \in (-1/e, 0)$, which indicates that there are two roots (Corless et al., 1996). We denote the root greater than -1 as $W_0(d)$ and the root less than -1 as $W_{-1}(d)$. It is obvious that $W_0(d) = -\alpha$. However, $W_0(d)$ is not the desired solution, since for $b > 0$, $x = -\alpha$ implies $\gamma = 0$, but we assume $\gamma > 0$. Hence, $W_{-1}(d)$ is the desired root, which gives the corresponding $\gamma = -\frac{1}{b}(W_{-1}(d) + \alpha) > 0$. Thus, we obtain

$$\gamma_{max} = -\frac{1}{b} \left(W_{-1} \left(-\frac{1}{1 + \frac{b^2}{V_{\min}}} \cdot e^{-\frac{1}{1 + \frac{b^2}{V_{\min}}}} \right) + \frac{1}{1 + \frac{b^2}{V_{\min}}} \right).$$

We construct the grid by taking $\gamma_i = c_2^{i-1} \gamma_{\min}$ for $i \in \{1, \dots, k_\gamma\}$, where $k_\gamma = \lceil \ln(\gamma_{max}/\gamma_{\min}) / \ln c_2 \rceil$. In the experiments we took $c_2 = 1.05$, and $\delta_1 = \delta_2 = \delta/2$.

Minimization of the bound

The adjustment of the bound in Theorem 2.13 to our hypothesis space construction, as described above, is:

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \mu)^2} \left(\mathbb{E}_{\rho^2}[\hat{L}_\mu(h, h', S_h \cap S_{h'})] + \frac{2 \text{KL}(\rho \parallel \pi) + \ln \frac{2k}{\delta}}{\gamma n} \right. \\ \left. + \frac{\phi(\gamma K_\mu)}{\gamma K_\mu^2} \left(\frac{\mathbb{E}_{\rho^2}[\hat{V}_\mu(h, h', S_h \cap S_{h'})]}{1 - \frac{\lambda n}{2(n-1)}} + \frac{K_\mu^2 (2 \text{KL}(\rho \parallel \pi) + \ln \frac{2k}{\delta})}{n\lambda \left(1 - \frac{\lambda n}{2(n-1)}\right)} \right) \right),$$

where $n = \min_{h, h'} |S_h \cap S_{h'}|$ and $k = k_\lambda k_\gamma$. We minimize the bound without considering k_γ and k_λ since we define the grid without taking them into consideration. However, we put back k_γ and k_λ when computing the generalization bound. Thus, when doing the optimization we take $k = 1$, but when we compute the bound we take the proper $k = k_\lambda k_\gamma$.

Algorithm 2: Minimization of the bound in Theorem 2.13

Input: n , grid of μ and losses $\mathbb{1}(h(X_i) \neq Y_i)$ for all $(X_i, Y_i) \in S_h$ for all h
for μ selected by the binary search in the grid **do**
 Initialize: $\rho = \pi$
 Compute $\hat{L}_\mu(h, h', S_h \cap S_{h'})$ and $\hat{V}_\mu(h, h', S_h \cap S_{h'})$ for all h, h'
 while The improvement of the bound for a fixed μ is larger than 10^{-9} **do**
 Compute $\lambda_{\mu, \rho}^*$, the optimal λ given ρ and μ
 Compute $\gamma_{\mu, \rho}^*$, the optimal γ given ρ and μ
 Apply gradient descent to the bound w.r.t. ρ given μ , $\lambda_{\mu, \rho}^*$ and $\gamma_{\mu, \rho}^*$
 end while
 Proceed to the next μ in the grid proposed by the binary search
end for

Optimal λ given μ and ρ Given μ and ρ , λ can be computed in the same way as in the optimization of Theorem 2.12, since the optimization problem is the same, and get

$$\lambda_{\mu, \rho}^* = \frac{2(n-1)}{n} \left(\sqrt{\frac{2(n-1) \mathbb{E}_{\rho^2}[\hat{V}_\mu(h, h', S_h \cap S_{h'})]}{K_\mu^2 (2 \text{KL}(\rho \parallel \pi) + \ln \frac{2k}{\delta})} + 1 + 1} \right)^{-1}.$$

In our implementation at every optimization step we took the closest λ to the above value from the λ -grid.

Optimal γ given μ and ρ Given μ and ρ , the bound for the variance is obtained by plugging in the optimal $\lambda_{\mu,\rho}^*$ computed above. Let

$$U_{\mathbb{V}}(\rho, \mu) = \frac{\mathbb{E}_{\rho^2}[\hat{\mathbb{V}}_{\mu}(h, h', S_h \cap S_{h'})]}{1 - \frac{\lambda_{\mu,\rho}^* n}{2(n-1)}} + \frac{K_{\mu}^2 (2 \text{KL}(\rho \|\pi) + \ln \frac{2k}{\delta})}{n \lambda_{\mu,\rho}^* \left(1 - \frac{\lambda_{\mu,\rho}^* n}{2(n-1)}\right)}.$$

Then

$$\gamma_{\mu,\rho}^* = \frac{1}{K_{\mu}} \left(W_0 \left(\frac{1}{e} \left(\frac{K_{\mu}^2 (2 \text{KL}(\rho \|\pi) + \ln \frac{2k}{\delta})}{n U_{\mathbb{V}}(\rho, \mu)} - 1 \right) \right) + 1 \right),$$

where W_0 is the principal branch of the Lambert W function, which is defined as the inverse of the function $f(x) = xe^x$. In our implementation at every optimization step we took the closest γ to the above value from the γ -grid.

Gradient w.r.t. ρ given λ, γ , and μ Optimizing the bound w.r.t. ρ is equivalent to constrained optimization of $f(\rho) = \mathbb{E}_{\rho^2}[\hat{L}_{\mu}(h, h', S')] + a \mathbb{E}_{\rho^2}[\hat{\mathbb{V}}_{\mu}(h, h', S')] + 2b \text{KL}(\rho \|\pi)$, where

$$a = \frac{\phi(K_{\mu}\gamma)}{K_{\mu}^2 \gamma} \frac{1}{1 - \frac{n\lambda}{2(n-1)}}, \quad b = \frac{1}{\gamma n} + \frac{\phi(K_{\mu}\gamma)}{K_{\mu}^2 \gamma} \frac{K_{\mu}^2}{n\lambda(1 - \frac{n\lambda}{2(n-1)})},$$

and the constraint is that ρ must be a probability distribution. We optimize ρ in the same way as presented in Appendix 2.8.4.1. We use \hat{L}_{μ} to denote the matrix of empirical μ -tandem losses and $\hat{\mathbb{V}}_{\mu}$ to denote the matrix of empirical variance of the μ -tandem losses. Then, the gradient w.r.t. ρ is given by:

$$\begin{aligned} (\nabla f)_h &= 2 \left(\sum_{h'} \rho(h') (\hat{L}_{\mu}(h, h', S') + a \hat{\mathbb{V}}_{\mu}(h, h', S')) + b \left(1 + \ln \frac{\rho(h)}{\pi(h)} \right) \right), \\ \nabla f &= 2 \left(\hat{L}_{\mu} \rho + a \hat{\mathbb{V}}_{\mu} \rho + b \left(1 + \ln \frac{\rho}{\pi} \right) \right). \end{aligned}$$

We applied gradient descent in the same way as presented in Appendix 2.8.4.1.

2.8.5 Experiments

2.8.5.1 Data sets

As mentioned, we considered data sets from the UCI and LibSVM repositories (Dua and Graff, 2019; Chang and Lin, 2011), as well as Fashion-MNIST (Fashion)

Table 2.1: Data set overview. c_{\min} and c_{\max} denote the minimum and maximum class frequency.

Data set	N	d	c	c_{\min}	c_{\max}	Source
Adult	32561	123	2	0.2408	0.7592	LIBSVM (a1a)
Cod-RNA	59535	8	2	0.3333	0.6667	LIBSVM
Connect-4	67557	126	3	0.0955	0.6583	LIBSVM
Fashion	70000	784	10	0.1000	0.1000	Zalando Research
Letter	20000	16	26	0.0367	0.0406	UCI
MNIST	70000	780	10	0.0902	0.1125	LIBSVM
Mushroom	8124	22	2	0.4820	0.5180	LIBSVM
Pendigits	10992	16	10	0.0960	0.1041	LIBSVM
Phishing	11055	68	2	0.4431	0.5569	LIBSVM
Protein	24387	357	3	0.2153	0.4638	LIBSVM
SVMGuide1	3089	4	2	0.3525	0.6475	LIBSVM
SatImage	6435	36	6	0.0973	0.2382	LIBSVM
Sensorless	58509	48	11	0.0909	0.0909	LIBSVM
Shuttle	58000	9	7	0.0002	0.7860	LIBSVM
Splice	3175	60	2	0.4809	0.5191	LIBSVM
USPS	9298	256	10	0.0761	0.1670	LIBSVM
w1a	49749	300	2	0.0297	0.9703	LIBSVM

from Zalando Research³. We used data sets with size $3000 \leq N \leq 70000$ and dimension $d \leq 1000$. These relatively large data sets were chosen in order to provide meaningful bounds in the standard bagging setting, where individual trees are trained on $n = 0.8N$ randomly subsampled points with replacement and the size of the overlap of out-of-bag sets is roughly $n/9$. An overview of the data sets is given in Table 2.1.

For all experiments, we removed patterns with missing entries and made a stratified split of the data set. For data sets with a training and a test set (SVMGuide1, Splice, Adult, w1a, MNIST, Shuttle, Pendigits, Protein, SatImage, USPS) we combined the training and test sets and shuffled the entire set before splitting.

³<https://research.zalando.com/welcome/mission/research-projects/fashion-mnist/>

Table 2.2: Numerical values of the test loss obtained by the RFs with optimized weighting. The smallest loss is highlighted in **bold**, while the smallest optimized loss is underlined.

Data set	$L(MV_u)$	$L(MV_{\rho_\lambda})$	$L(MV_{\rho_{\text{TND}}})$	$L(MV_{\rho_{\text{CCTND}}})$	$L(MV_{\rho_{\text{CCPBB}}})$
SVMGuide1	0.0284 (0.0037)	0.0372 (0.0066)	0.0287 (0.0035)	<u>0.0286 (0.0036)</u>	0.0287 (0.0039)
Phishing	0.0292 (0.004)	0.0371 (0.0073)	0.0292 (0.0036)	0.0292 (0.0036)	0.0292 (0.004)
Mushroom	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)
Splice	0.0299 (0.009)	0.1087 (0.021)	0.0306 (0.0099)	0.0309 (0.0092)	<u>0.0302 (0.01)</u>
w1a	0.0108 (0.0007)	0.016 (0.0025)	0.0108 (0.0006)	0.0107 (0.0006)	0.0108 (0.0006)
Cod-RNA	0.0402 (0.0013)	0.0712 (0.0064)	0.0395 (0.0014)	0.0395 (0.0014)	0.0395 (0.0015)
Adult	0.1693 (0.0027)	0.1942 (0.0151)	<u>0.1698 (0.0031)</u>	0.1701 (0.003)	0.1698 (0.0031)
Connect-4	0.1706 (0.0023)	0.2803 (0.0165)	<u>0.1699 (0.002)</u>	0.1705 (0.0024)	0.1695 (0.0019)
Shuttle	0.0002 (0.0001)	0.0003 (0.0002)	0.0002 (0.0001)	0.0002 (0.0001)	0.0002 (0.0001)
Pendigits	0.0096 (0.0023)	0.0452 (0.0124)	0.0092 (0.0022)	0.0093 (0.0021)	0.0092 (0.0025)
Letter	0.0378 (0.0036)	0.1408 (0.0356)	<u>0.0398 (0.0041)</u>	0.0402 (0.0042)	<u>0.0383 (0.0034)</u>
SatImage	0.0828 (0.0068)	0.1321 (0.0268)	0.0835 (0.0061)	0.0839 (0.0062)	<u>0.0832 (0.006)</u>
Sensorless	0.0014 (0.0004)	0.0138 (0.0019)	0.0012 (0.0003)	0.0012 (0.0003)	0.0012 (0.0003)
USPS	0.0394 (0.0043)	0.1325 (0.0251)	<u>0.0401 (0.0055)</u>	0.0405 (0.0052)	0.0404 (0.005)
MNIST	0.0316 (0.0017)	0.16 (0.0352)	<u>0.0323 (0.0017)</u>	0.0324 (0.0017)	0.0317 (0.0014)
Fashion	0.1175 (0.0018)	0.2122 (0.0299)	0.1192 (0.0022)	0.1197 (0.0022)	<u>0.1178 (0.0021)</u>

2.8.5.2 Optimized weighted random forest

Experimental Setting This section describes in detail the settings and the results of the empirical evaluation using random forest (RF) majority vote classifiers.

We construct the ensemble from decision trees available in *scikit-learn*. For each data set, an ensemble of 100 trees is trained using bagging (as described in Section 2.6). For each tree, the Gini criterion is used for splitting and \sqrt{d} features are considered in each split.

We compare the RF using the default uniform weighting ρ_u and the optimized weighting obtained by FO (Thiemann et al., 2016), TND (Masegosa et al., 2020), CCTND (Theorem 2.11) and CCPBB (Theorem 2.13). Optimization is based on the out-of-bag sets (see Section 2.6). For each optimized RF, we also compute the optimized bound.

Numerical Results This section lists the numerical results for the empirical evaluation using RF. Table 2.2 provides the numerical values of the test loss obtained by the RFs with uniform weighting and with weighting optimized by FO, TND, CCTND and CCPBB; a visual presentation is given in Figure 2.2a. As observed

Table 2.3: Numerical values of the bounds for the RFs with optimized weighting. The tightest bound is highlighted in **bold**, while the tightest second-order bound is underlined.

Data set	FO(ρ_λ)	TND(ρ_{TND})	CCTND(ρ_{CCTND})	CCPBB(ρ_{CCPBB})
SVMGuidel	0.1079 (0.0079)	0.1836 (0.0062)	0.1853 (0.0059)	0.2806 (0.0071)
Phishing	0.1189 (0.0035)	<u>0.1642 (0.0043)</u>	0.1674 (0.0042)	0.2336 (0.005)
Mushroom	0.0068 (0.0001)	<u>0.0353 (0.0002)</u>	0.0388 (0.0002)	0.1121 (0.0006)
Splice	0.3245 (0.0218)	<u>0.4077 (0.0062)</u>	0.4247 (0.0065)	0.6562 (0.0056)
w1a	0.0424 (0.0015)	<u>0.0633 (0.0009)</u>	0.0642 (0.0009)	0.0805 (0.0011)
Cod-RNA	0.1629 (0.0018)	<u>0.1663 (0.0014)</u>	0.1698 (0.0014)	0.19 (0.0018)
Adult	0.4388 (0.0042)	<u>0.5701 (0.0051)</u>	<u>0.5508 (0.004)</u>	0.5976 (0.0042)
Connect-4	0.5978 (0.0067)	0.6831 (0.0039)	<u>0.6758 (0.0036)</u>	0.7112 (0.0038)
Shuttle	0.0026 (0.0002)	<u>0.0078 (0.0002)</u>	<u>0.0083 (0.0002)</u>	0.018 (0.0003)
Pendigits	0.142 (0.0035)	<u>0.1445 (0.0026)</u>	0.1504 (0.0042)	0.2155 (0.003)
Letter	0.3858 (0.0067)	<u>0.4504 (0.0032)</u>	0.4513 (0.003)	0.5134 (0.0039)
SatImage	0.3762 (0.0075)	0.4902 (0.0079)	<u>0.4851 (0.007)</u>	0.6158 (0.0083)
Sensorless	0.0348 (0.0031)	0.0257 (0.0006)	<u>0.0265 (0.0006)</u>	0.0376 (0.0007)
USPS	0.3394 (0.0065)	<u>0.4059 (0.0048)</u>	0.4097 (0.0044)	0.5086 (0.0042)
MNIST	0.3795 (0.0031)	0.3537 (0.0014)	0.3598 (0.0014)	0.3853 (0.0014)
Fashion	0.4806 (0.003)	<u>0.5436 (0.0023)</u>	<u>0.5408 (0.0021)</u>	0.5728 (0.0021)

by Masegosa et al. (2020), optimization using FO leads to overfitting, while the second-order bounds does not significantly degrade the performance. Among the second-order bounds, optimizing using CCPBB produces the best classifier in most cases.

Table 2.3 provides the numerical values of the optimized bounds; a visual presentation is given in Figure 2.2b. Table 2.4 provides the Gibbs loss and tandem loss using the optimized ρ . The optimal μ found is reported for CCTND and CCPBB as well.

2.8.5.3 Ensemble of multiple heterogeneous classifiers

Experimental Setting This section describes in detail the settings and the results of the experimental evaluation using an ensemble of multiple heterogeneous classifiers.

The ensemble is defined by a set of standard classifiers available in *scikit-learn*:

- **Linear Discriminant Analysis**, with default parameters, which includes a singular value decomposition solver.

Table 2.4: Numerical values for Gibbs loss, tandem loss and optimized μ for the RFs with optimized weighting. We use $\mathbb{E}_\rho[L]$ and $\mathbb{E}_{\rho^2}[L]$ as short-hands for the Gibbs and the tandem loss respectively.

Data set	FO		TND		CCTND			CCPBB		
	$\mathbb{E}_\rho[L]$	$\mathbb{E}_{\rho^2}[L]$	$\mathbb{E}_\rho[L]$	$\mathbb{E}_{\rho^2}[L]$	$\mathbb{E}_\rho[L]$	$\mathbb{E}_{\rho^2}[L]$	μ	$\mathbb{E}_\rho[L]$	$\mathbb{E}_{\rho^2}[L]$	μ
SVMGuidel	0.0325	0.0217	0.0406	0.0185	0.0403	0.0184	-0.0527	0.0413	0.0194	-0.0258
Phishing	0.041	0.0255	0.0486	0.0197	0.0484	0.0196	-0.0295	0.049	0.0202	-0.0125
Mushroom	0.0	0.0	0.0002	0.0	0.0002	0.0	-0.0317	0.0002	0.0	-0.01
Splice	0.1068	0.0903	0.1564	0.0424	0.1522	0.0415	-0.057	0.16	0.044	0.0045
w1a	0.0156	0.0123	0.0179	0.0091	0.0179	0.009	-0.0111	0.018	0.0092	-0.0065
Cod-RNA	0.0712	0.0602	0.0802	0.0314	0.0803	0.0314	-0.0178	0.0815	0.0318	0.0102
Adult	0.1995	0.1474	0.2061	0.1184	0.2056	0.1182	-0.1216	0.2068	0.1194	-0.0918
Connect-4	0.2824	0.2564	0.2953	0.1523	0.2943	0.1521	-0.0959	0.2974	0.1535	-0.0615
Shuttle	0.0003	0.0001	0.0006	0.0002	0.0006	0.0002	-0.0044	0.0006	0.0002	0.0
Pendigits	0.0502	0.0346	0.061	0.0163	0.0609	0.0163	-0.0099	0.0614	0.0166	0.0092
Letter	0.1685	0.1249	0.1803	0.0851	0.1797	0.0849	-0.0501	0.1816	0.0861	-0.0228
SatImage	0.1478	0.0968	0.1612	0.0746	0.1602	0.0741	-0.1104	0.1617	0.0755	-0.0535
Sensorless	0.0125	0.0113	0.0192	0.0027	0.0192	0.0027	0.0008	0.0195	0.0027	0.01
USPS	0.1363	0.0989	0.1517	0.0644	0.1509	0.0641	-0.053	0.1522	0.065	-0.0173
MNIST	0.1763	0.1286	0.1837	0.075	0.1835	0.075	0.0281	0.185	0.0756	0.037
Fashion	0.2256	0.1715	0.2325	0.1196	0.2322	0.1195	-0.0577	0.2334	0.1203	-0.0382

- Three versions of **k-Nearest Neighbors**: (i) $k=3$ and uniform weights (i.e., all points in each neighborhood are weighted equally); (ii) $k=5$ and uniform weights; and (iii) $k=5$ where points are weighted by the inverse of their distance. In all cases, it is employed the Euclidean distance.
- **Decision Tree**, with default parameters, which includes Gini criterion for splitting and no maximum depth.
- **Logistic Regression**, with default parameters, which includes L2 penalization.
- **Gaussian Naive Bayes**, with default parameters.

We included three versions of the kNN classifier to test if our bounds could deal with a heterogeneous set of classifiers where some of them are expected to provide highly correlated errors while others are expected to provide much less correlated errors.

Each of the seven classifiers of the ensemble was learned from a bootstrap sample of the training data set. We did it in the way to be able to compute and optimize our bounds with the out-of-bag-samples as described in Section 2.6.

Numerical Results This section lists the numerical results for the empirical evaluation using ensembles of multiple heterogeneous classifiers.

Table 2.5: Numerical values of the test loss obtained by ensembles of multiple heterogeneous classifiers with optimized weighting. The smallest loss is highlighted in **bold**, while the smallest optimized loss is underlined.

Data set	$L(MV_u)$	$L(MV_{\rho_\lambda})$	$L(MV_{\rho_{TND}})$	$L(MV_{\rho_{CCTND}})$	$L(MV_{\rho_{CCPBB}})$
SVMGuide1	0.0284 (0.0037)	0.0372 (0.0066)	0.0287 (0.0035)	<u>0.0286 (0.0036)</u>	0.0287 (0.0039)
Phishing	0.0292 (0.004)	0.0371 (0.0073)	0.0292 (0.0036)	0.0292 (0.0036)	0.0292 (0.004)
Mushroom	0.0 (0.0)	<u>0.0 (0.0)</u>	<u>0.0 (0.0)</u>	<u>0.0 (0.0)</u>	<u>0.0 (0.0)</u>
Splice	0.0299 (0.009)	0.1087 (0.021)	0.0306 (0.0099)	0.0309 (0.0092)	<u>0.0302 (0.01)</u>
w1a	0.0108 (0.0007)	0.016 (0.0025)	0.0108 (0.0006)	0.0107 (0.0006)	0.0108 (0.0006)
Cod-RNA	0.0402 (0.0013)	0.0712 (0.0064)	0.0395 (0.0014)	0.0395 (0.0014)	0.0395 (0.0015)
Adult	0.1693 (0.0027)	0.1942 (0.0151)	<u>0.1698 (0.0031)</u>	0.1701 (0.003)	0.1698 (0.0031)
Connect-4	0.1706 (0.0023)	0.2803 (0.0165)	<u>0.1699 (0.002)</u>	0.1705 (0.0024)	0.1695 (0.0019)
Shuttle	0.0002 (0.0001)	0.0003 (0.0002)	0.0002 (0.0001)	0.0002 (0.0001)	0.0002 (0.0001)
Pendigits	0.0096 (0.0023)	0.0452 (0.0124)	0.0092 (0.0022)	0.0093 (0.0021)	0.0092 (0.0025)
Letter	0.0378 (0.0036)	0.1408 (0.0356)	<u>0.0398 (0.0041)</u>	0.0402 (0.0042)	<u>0.0383 (0.0034)</u>
SatImage	0.0828 (0.0068)	0.1321 (0.0268)	0.0835 (0.0061)	0.0839 (0.0062)	<u>0.0832 (0.006)</u>
Sensorless	0.0014 (0.0004)	0.0138 (0.0019)	0.0012 (0.0003)	0.0012 (0.0003)	0.0012 (0.0003)
USPS	0.0394 (0.0043)	0.1325 (0.0251)	<u>0.0401 (0.0055)</u>	0.0405 (0.0052)	0.0404 (0.005)
MNIST	0.0316 (0.0017)	0.16 (0.0352)	<u>0.0323 (0.0017)</u>	0.0324 (0.0017)	0.0317 (0.0014)
Fashion	0.1175 (0.0018)	0.2122 (0.0299)	0.1192 (0.0022)	0.1197 (0.0022)	<u>0.1178 (0.0021)</u>

Table 2.5 provides the numerical values of the test loss obtained by these ensembles with uniform weighting and with weighting optimized by FO, TND, CCTND and CCPBB; a visual presentation is given in Figure 2.2c. In this case, uniform voting is not a competitive weighting scheme. The second-order bounds perform much better than uniform weighting and than the weights computed according to the first-order bound. There is not any clear winner among the second-order bounds.

Table 2.6 provides the numerical values of the optimized bounds; a visual presentation is given in Figure 2.2d. Among the second-order bounds, the CCTND bound is often tighter in this setting.

Table 2.7 provides the recorded Gibbs loss and tandem loss using the optimized ρ . The optimal μ found is reported for CCTND and CCPBB as well.

Table 2.6: Numerical values of the bounds for ensembles of multiple heterogeneous classifiers with optimized weighting. The tightest bound is highlighted in **bold**, while the tightest second-order bound is underlined.

Data set	FO(ρ_λ)	TND(ρ_{TND})	CCTND(ρ_{CCTND})	CCPBB(ρ_{CCPBB})
SVMGuide1	0.1079 (0.0079)	0.1836 (0.0062)	0.1853 (0.0059)	0.2806 (0.0071)
Phishing	0.1189 (0.0035)	<u>0.1642 (0.0043)</u>	0.1674 (0.0042)	0.2336 (0.005)
Mushroom	0.0068 (0.0001)	<u>0.0353 (0.0002)</u>	0.0388 (0.0002)	0.1121 (0.0006)
Splice	0.3245 (0.0218)	<u>0.4077 (0.0062)</u>	0.4247 (0.0065)	0.6562 (0.0056)
w1a	0.0424 (0.0015)	<u>0.0633 (0.0009)</u>	0.0642 (0.0009)	0.0805 (0.0011)
Cod-RNA	0.1629 (0.0018)	<u>0.1663 (0.0014)</u>	0.1698 (0.0014)	0.19 (0.0018)
Adult	0.4388 (0.0042)	0.5701 (0.0051)	<u>0.5508 (0.004)</u>	0.5976 (0.0042)
Connect-4	0.5978 (0.0067)	0.6831 (0.0039)	<u>0.6758 (0.0036)</u>	0.7112 (0.0038)
Shuttle	0.0026 (0.0002)	<u>0.0078 (0.0002)</u>	<u>0.0083 (0.0002)</u>	0.018 (0.0003)
Pendigits	0.142 (0.0035)	<u>0.1445 (0.0026)</u>	0.1504 (0.0042)	0.2155 (0.003)
Letter	0.3858 (0.0067)	<u>0.4504 (0.0032)</u>	0.4513 (0.003)	0.5134 (0.0039)
SatImage	0.3762 (0.0075)	<u>0.4902 (0.0079)</u>	<u>0.4851 (0.007)</u>	0.6158 (0.0083)
Sensorless	0.0348 (0.0031)	0.0257 (0.0006)	<u>0.0265 (0.0006)</u>	0.0376 (0.0007)
USPS	0.3394 (0.0065)	<u>0.4059 (0.0048)</u>	0.4097 (0.0044)	0.5086 (0.0042)
MNIST	0.3795 (0.0031)	0.3537 (0.0014)	0.3598 (0.0014)	0.3853 (0.0014)
Fashion	0.4806 (0.003)	<u>0.5436 (0.0023)</u>	<u>0.5408 (0.0021)</u>	0.5728 (0.0021)

Table 2.7: Numerical values for Gibbs loss, tandem loss and optimized μ for the heterogeneous classifiers with optimized weighting. We use $\mathbb{E}_\rho[L]$ and $\mathbb{E}_{\rho^2}[L]$ as short-hands for the Gibbs loss and the tandem loss respectively.

Data set	FO		TND		CCTND			CCPBB		
	$\mathbb{E}_\rho[L]$	$\mathbb{E}_{\rho^2}[L]$	$\mathbb{E}_\rho[L]$	$\mathbb{E}_{\rho^2}[L]$	$\mathbb{E}_\rho[L]$	$\mathbb{E}_{\rho^2}[L]$	μ	$\mathbb{E}_\rho[L]$	$\mathbb{E}_{\rho^2}[L]$	μ
SVMGuide1	0.0325	0.0217	0.0406	0.0185	0.0403	0.0184	-0.0527	0.0413	0.0194	-0.0258
Phishing	0.041	0.0255	0.0486	0.0197	0.0484	0.0196	-0.0295	0.049	0.0202	-0.0125
Mushroom	0.0	0.0	0.0002	0.0	0.0002	0.0	-0.0317	0.0002	0.0	-0.01
Splice	0.1068	0.0903	0.1564	0.0424	0.1522	0.0415	-0.057	0.16	0.044	0.0045
w1a	0.0156	0.0123	0.0179	0.0091	0.0179	0.009	-0.0111	0.018	0.0092	-0.0065
Cod-RNA	0.0712	0.0602	0.0802	0.0314	0.0803	0.0314	-0.0178	0.0815	0.0318	0.0102
Adult	0.1995	0.1474	0.2061	0.1184	0.2056	0.1182	-0.1216	0.2068	0.1194	-0.0918
Connect-4	0.2824	0.2564	0.2953	0.1523	0.2943	0.1521	-0.0959	0.2974	0.1535	-0.0615
Shuttle	0.0003	0.0001	0.0006	0.0002	0.0006	0.0002	-0.0044	0.0006	0.0002	0.0
Pendigits	0.0502	0.0346	0.061	0.0163	0.0609	0.0163	-0.0099	0.0614	0.0166	0.0092
Letter	0.1685	0.1249	0.1803	0.0851	0.1797	0.0849	-0.0501	0.1816	0.0861	-0.0228
SatImage	0.1478	0.0968	0.1612	0.0746	0.1602	0.0741	-0.1104	0.1617	0.0755	-0.0535
Sensorless	0.0125	0.0113	0.0192	0.0027	0.0192	0.0027	0.0008	0.0195	0.0027	0.01
USPS	0.1363	0.0989	0.1517	0.0644	0.1509	0.0641	-0.053	0.1522	0.065	-0.0173
MNIST	0.1763	0.1286	0.1837	0.075	0.1835	0.075	0.0281	0.185	0.0756	0.037
Fashion	0.2256	0.1715	0.2325	0.1196	0.2322	0.1195	-0.0577	0.2334	0.1203	-0.0382

Chapter 3

Split-kl and PAC-Bayes-split-kl Inequalities for Ternary Random Variables

The work presented in this chapter is based on a paper that has been published as:

Yi-Shan Wu and Yevgeny Seldin. Split-kl and pac-bayes-split-kl inequalities for ternary random variables. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.

Abstract

We present a new concentration of measure inequality for sums of independent bounded random variables, which we name a split-kl inequality. The inequality is particularly well-suited for ternary random variables, which naturally show up in a variety of problems, including analysis of excess losses in classification, analysis of weighted majority votes, and learning with abstention. We demonstrate that for ternary random variables the inequality is simultaneously competitive with the kl inequality, the Empirical Bernstein inequality, and the Unexpected Bernstein inequality, and in certain regimes outperforms all of them. It resolves an open question by Tolstikhin and Seldin (2013) and Mhammedi et al. (2019) on how to match simultaneously the combinatorial power of the kl inequality when the distribution happens to be close to binary and the power of Bernstein inequalities to exploit low variance when the probability mass is concentrated on the middle value. We also derive a PAC-Bayes-split-kl inequality and compare it with the PAC-Bayes-kl, PAC-Bayes-Empirical-Bennett, and PAC-Bayes-Unexpected-Bernstein inequalities in an analysis of excess losses and in an analysis of a weighted majority vote for several UCI datasets. Last but not least, our study provides the first direct comparison of the Empirical Bernstein and Unexpected Bernstein inequalities and their PAC-Bayes extensions.

3.1 Introduction

Concentration of measure inequalities for sums of independent random variables are the most fundamental analysis tools in statistics and many other domains (Boucheron et al., 2013). Their history stretches almost a century back, and inequalities such as Hoeffding’s (Hoeffding, 1963) and Bernstein’s (Bernstein, 1946) are the main work horses of learning theory.

For binary random variables, one of the tightest concentration of measure inequalities is the kl inequality (Maurer, 2004; Langford, 2005; Foong et al., 2021, 2022), which is based on combinatorial properties of a sum of n independent random variables.¹ However, while being extremely tight for binary random variables and applicable to any bounded random variables, the kl inequality is not necessarily a good choice for sums of bounded random variables that can take more than two values. In the latter

¹The Binomial tail bound is slightly tighter, but it does not extend to the PAC-Bayes setting (Langford, 2005). Our split-kl approach can be directly applied to obtain a “split-Binomial-tail” inequality.

case, the Empirical Bernstein (Mnih et al., 2008; Audibert et al., 2009; Maurer and Pontil, 2009) and the Unexpected Bernstein (Cesa-Bianchi et al., 2007; Mhammedi et al., 2019) inequalities can be significantly tighter due to their ability to exploit low variance, as shown by Tolstikhin and Seldin (2013). However, the Empirical and Unexpected Bernstein inequalities are loose for binary random variables (Tolstikhin and Seldin, 2013).

The challenge of exploiting low variance and, at the same time, matching the tightness of the kl inequality if a distribution happens to be close to binary, was faced by multiple prior works (Tolstikhin and Seldin, 2013; Mhammedi et al., 2019; Wu et al., 2021), but remained an open question. We resolve this question for the case of ternary random variables. Such random variables appear in a variety of applications, and we illustrate two of them. One is a study of excess losses, which are differences between the zero-one losses of a prediction rule h and a reference prediction rule h^* , $Z = \ell(h(X), Y) - \ell(h^*(X), Y) \in \{-1, 0, 1\}$. Mhammedi et al. (2019) have applied the PAC-Bayes-Unexpected-Bernstein bound to excess losses in order to improve generalization bounds for classification. Another example of ternary random variables is the tandem loss with an offset, defined by $\ell_\alpha(h(X), h'(X), Y) = (\ell(h(X), Y) - \alpha)(\ell(h'(X), Y) - \alpha) \in \{\alpha^2, -\alpha(1 - \alpha), (1 - \alpha)^2\}$. Wu et al. (2021) have applied the PAC-Bayes-Empirical-Bennett inequality to the tandem loss with an offset to obtain a generalization bound for the weighted majority vote. Yet another potential application, which we leave for future work, is learning with abstention (Cortes et al., 2018; Thulasidasan et al., 2019).

We present the split-kl inequality, which simultaneously matches the tightness of the Empirical/Unexpected Bernstein and the kl, and outperforms both for certain distributions. It works for sums of any bounded random variables Z_1, \dots, Z_n , not only the ternary ones, but it is best suited for ternary random variables, for which it is almost tight (in the same sense, as the kl is tight for binary random variables). The idea behind the split-kl inequality is to write a random variable Z as $Z = \mu + Z^+ - Z^-$, where μ is a constant, $Z^+ = \max\{0, Z - \mu\}$, and $Z^- = \max\{0, \mu - Z\}$. Then $\mathbb{E}[Z] = \mu + \mathbb{E}[Z^+] - \mathbb{E}[Z^-]$ and, given an i.i.d. sample Z_1, \dots, Z_n , we can bound the distance between $\frac{1}{n} \sum_{i=1}^n Z_i$ and $\mathbb{E}[Z]$ by using kl upper and lower bounds on the distances between $\frac{1}{n} \sum_{i=1}^n Z_i^+$ and $\mathbb{E}[Z^+]$, and $\frac{1}{n} \sum_{i=1}^n Z_i^-$ and $\mathbb{E}[Z^-]$, respectively. For ternary random variables $Z \in \{a, b, c\}$ with $a \leq b \leq c$, the best split is to take $\mu = b$, then both Z^+ and Z^- are binary and the kl upper and lower bounds for their rescaled versions are tight and, therefore, the split-kl inequality for Z is also tight. Thus, this approach provides the best of both worlds: the combinatorial tightness of the kl bound and exploitation of low variance when the probability mass on the

middle value happens to be large, as in Empirical Bernstein inequalities. We further elevate the idea to the PAC-Bayes domain and derive a PAC-Bayes-split-kl inequality.

We present an extensive set of experiments, where we first compare the kl, Empirical Bernstein, Unexpected Bernstein, and split-kl inequalities applied to (individual) sums of independent random variables in simulated data, and then compare the PAC-Bayes-kl, PAC-Bayes-Unexpected-Bernstein, PAC-Bayes-split-kl, and, in some of the setups, PAC-Bayes-Empirical-Bennett, for several prediction models on several UCI datasets. In particular, we evaluate the bounds in the linear classification setup studied by Mhammedi et al. (2019) and in the weighted majority prediction setup studied by Wu et al. (2021). To the best of our knowledge, this is also the first time when the Empirical Bernstein and the Unexpected Bernstein inequalities are directly compared, with and without the PAC-Bayesian extension. In Appendix 3.6.1.2 we also show that an inequality introduced by Cesa-Bianchi et al. (2007) yields a relaxation of the Unexpected Bernstein inequality by Mhammedi et al. (2019).

3.2 Concentration of Measure Inequalities for Sums of Independent Random Variables

We start with the most basic question in probability theory and statistics: how far can an average of an i.i.d. sample Z_1, \dots, Z_n deviate from its expectation? We cite the major existing inequalities, the kl, Empirical Bernstein, and Unexpected Bernstein, then derive the new split-kl inequality, and then provide a numerical comparison.

3.2.1 Background

We use $\text{KL}(\rho||\pi)$ to denote the Kullback-Leibler divergence between two probability distributions, ρ and π (Cover and Thomas, 2006). We further use $\text{kl}(p||q)$ as a shorthand for the Kullback-Leibler divergence between two Bernoulli distributions with biases p and q , namely $\text{kl}(p||q) = \text{KL}((1-p, p)|| (1-q, q))$. For $\hat{p} \in [0, 1]$ and $\varepsilon \geq 0$ we define the upper and lower inverse of kl, respectively, as $\text{kl}^{-1,+}(\hat{p}, \varepsilon) := \max \{p : p \in [0, 1] \text{ and } \text{kl}(\hat{p}||p) \leq \varepsilon\}$ and $\text{kl}^{-1,-}(\hat{p}, \varepsilon) := \min \{p : p \in [0, 1] \text{ and } \text{kl}(\hat{p}||p) \leq \varepsilon\}$.

The first inequality that we cite is the kl inequality.

Theorem 3.1 (kl Inequality (Langford, 2005; Foong et al., 2021, 2022)). *Let Z_1, \dots, Z_n be i.i.d. random variables bounded in the $[0, 1]$ interval and with $\mathbb{E}[Z_i] = p$*

for all i . Let $\hat{p} = \frac{1}{n} \sum_{i=1}^n Z_i$ be their empirical mean. Then, for any $\delta \in (0, 1)$:

$$\mathbb{P}\left(\text{kl}(\hat{p}||p) \geq \frac{\ln \frac{1}{\delta}}{n}\right) \leq \delta$$

and, by inversion of the kl,

$$\mathbb{P}\left(p \geq \text{kl}^{-1,+}\left(\hat{p}, \frac{1}{n} \ln \frac{1}{\delta}\right)\right) \leq \delta, \quad (3.1)$$

$$\mathbb{P}\left(p \leq \text{kl}^{-1,-}\left(\hat{p}, \frac{1}{n} \ln \frac{1}{\delta}\right)\right) \leq \delta. \quad (3.2)$$

We note that the PAC-Bayes-kl inequality (Theorem 3.5 below) is based on the inequality $\mathbb{E}\left[e^{n \text{kl}(\hat{p}||p)}\right] \leq 2\sqrt{n}$ (Maurer, 2004), which gives $\mathbb{P}\left(\text{kl}(\hat{p}||p) \geq \frac{\ln \frac{2\sqrt{n}}{\delta}}{n}\right) \leq \delta$. Foong et al. (2021, 2022) reduce the logarithmic factor down to $\ln \frac{1}{\delta}$ by basing the proof on Chernoff's inequality, but this proof technique cannot be combined with PAC-Bayes. Therefore, when we move on to PAC-Bayes we pay the extra $\ln 2\sqrt{n}$ factor in the bounds. It is a long-standing open question whether this factor can be reduced in the PAC-Bayesian setting (Foong et al., 2021).

Next we cite two versions of the Empirical Bernstein inequality.

Theorem 3.2 (Empirical Bernstein Inequality (Maurer and Pontil, 2009)). *Let Z_1, \dots, Z_n be i.i.d. random variables bounded in a $[a, b]$ interval for some $a, b \in \mathbb{R}$, and with $\mathbb{E}[Z_i] = p$ for all i . Let $\hat{p} = \frac{1}{n} \sum_{i=1}^n Z_i$ be the empirical mean and let $\hat{\sigma} = \frac{1}{n-1} \sum_{i=1}^n (Z_i - \hat{p})^2$ be the empirical variance. Then for any $\delta \in (0, 1)$:*

$$\mathbb{P}\left(p \geq \hat{p} + \sqrt{\frac{2\hat{\sigma} \ln \frac{2}{\delta}}{n}} + \frac{7(b-a) \ln \frac{2}{\delta}}{3(n-1)}\right) \leq \delta. \quad (3.3)$$

Theorem 3.3 (Unexpected Bernstein Inequality (Fan et al., 2015; Mhammedi et al., 2019)). *Let Z_1, \dots, Z_n be i.i.d. random variables bounded from above by b for some $b > 0$, and with $\mathbb{E}[Z_i] = p$ for all i . Let $\hat{p} = \frac{1}{n} \sum_{i=1}^n Z_i$ be the empirical mean and let $\hat{\sigma} = \frac{1}{n} \sum_{i=1}^n Z_i^2$ be the empirical mean of the second moments. Let $\psi(u) := u - \ln(1+u)$ for $u > -1$. Then, for any $\gamma \in (0, 1/b)$ and any $\delta \in (0, 1)$:*

$$\mathbb{P}\left(p \geq \hat{p} + \frac{\psi(-\gamma b)}{\gamma b^2} \hat{\sigma} + \frac{\ln \frac{1}{\delta}}{\gamma n}\right) \leq \delta. \quad (3.4)$$

To facilitate a comparison with other bounds, Theorem 3.3 provides a slightly different form of the Unexpected Bernstein inequality than the one used by Mhammedi et al. (2019). We provide a proof of the theorem in Appendix 3.6.1.1, which is based on the Unexpected Bernstein Lemma (Fan et al., 2015). We note that an inequality proposed by Cesa-Bianchi et al. (2007) can be used to derive a relaxed version of the Unexpected Bernstein inequality, as discussed in Appendix 3.6.1.2.

3.2.2 The Split-kl Inequality

Let Z be a random variable bounded in a $[a, b]$ interval for some $a, b \in \mathbb{R}$ and let $\mu \in [a, b]$ be a constant. We decompose $Z = \mu + Z^+ - Z^-$, where $Z^+ = \max(0, Z - \mu)$ and $Z^- = \max(0, \mu - Z)$. Let $p = \mathbb{E}[Z]$, $p^+ = \mathbb{E}[Z^+]$, and $p^- = \mathbb{E}[Z^-]$. For an i.i.d. sample Z_1, \dots, Z_n let $\hat{p}^+ = \frac{1}{n} \sum_{i=1}^n Z_i^+$ and $\hat{p}^- = \frac{1}{n} \sum_{i=1}^n Z_i^-$.

With these definitions we present the split-kl inequality.

Theorem 3.4 (Split-kl inequality). *Let Z_1, \dots, Z_n be i.i.d. random variables in a $[a, b]$ interval for some $a, b \in \mathbb{R}$, then for any $\mu \in [a, b]$ and $\delta \in (0, 1)$:*

$$\mathbb{P}\left(p \geq \mu + (b - \mu) \text{kl}^{-1,+} \left(\frac{\hat{p}^+}{b - \mu}, \frac{1}{n} \ln \frac{2}{\delta} \right) - (\mu - a) \text{kl}^{-1,-} \left(\frac{\hat{p}^-}{\mu - a}, \frac{1}{n} \ln \frac{2}{\delta} \right)\right) \leq \delta. \quad (3.5)$$

Proof.

$$\begin{aligned} & \mathbb{P}\left(p \geq \mu + (b - \mu) \text{kl}^{-1,+} \left(\frac{\hat{p}^+}{b - \mu}, \frac{1}{n} \ln \frac{2}{\delta} \right) - (\mu - a) \text{kl}^{-1,-} \left(\frac{\hat{p}^-}{\mu - a}, \frac{1}{n} \ln \frac{2}{\delta} \right)\right) \\ & \leq \mathbb{P}\left(p^+ \geq (b - \mu) \text{kl}^{-1,+} \left(\frac{\hat{p}^+}{b - \mu}, \frac{1}{n} \ln \frac{2}{\delta} \right)\right) + \mathbb{P}\left(p^- \leq (\mu - a) \text{kl}^{-1,-} \left(\frac{\hat{p}^-}{\mu - a}, \frac{1}{n} \ln \frac{2}{\delta} \right)\right) \\ & \leq \delta, \end{aligned}$$

where the last inequality follows by application of the kl upper and lower bounds from Theorem 3.1 to the first and second terms in the middle line, respectively. \square

For ternary random variables the best choice is to take μ to be the middle value, then the resulting Z^+ and Z^- are binary and the corresponding kl upper and lower bounds on p^+ and p^- are tight, and the resulting split-kl bound is tight. The inequality can be applied to any bounded random variables, but same way as the kl inequality is not necessarily a good choice for bounded random variables, if the distribution is not binary, the split-kl is not necessarily a good choice if the distribution is not ternary.

3.2.3 Empirical Comparison

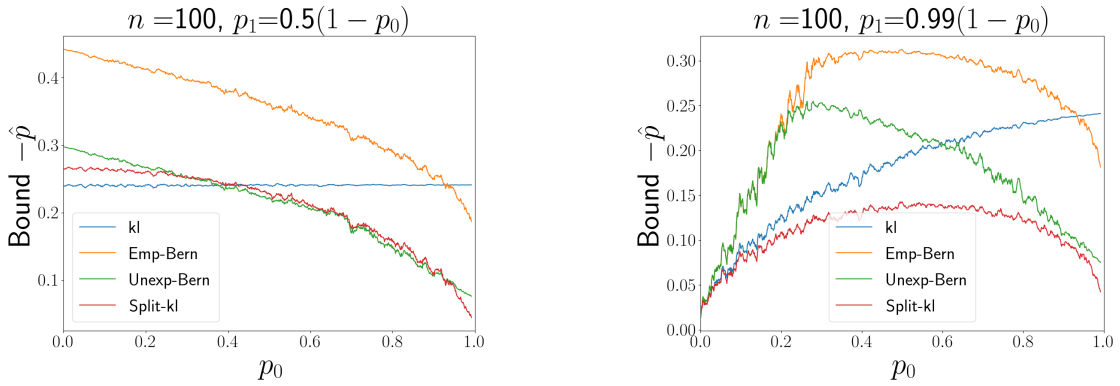
We present an empirical comparison of the tightness of the above four concentration inequalities: the kl, the Empirical Bernstein, the Unexpected Bernstein, and the split-kl. We take n i.i.d. samples Z_1, \dots, Z_n taking values in $\{-1, 0, 1\}$. The choice is motivated both by instructiveness of presentation and by subsequent applications to excess losses. We let $p_{-1} = \mathbb{P}(Z = -1)$, $p_0 = \mathbb{P}(Z = 0)$, and $p_1 = \mathbb{P}(Z = 1)$, where $p_{-1} + p_0 + p_1 = 1$. Then $p = \mathbb{E}[Z] = p_1 - p_{-1}$. We also let $\hat{p} = \frac{1}{n} \sum_{i=1}^n Z_i$.

In Figure 3.1 we plot the difference between the bounds on p given by the inequalities (3.1), (3.3), (3.4), and (3.5), and \hat{p} . Lower values in the plot correspond to tighter bounds. To compute the kl bound we first rescale the losses to the $[0, 1]$ interval, and then rescale the bound back to the $[-1, 1]$ interval. For the Empirical Bernstein bound we take $a = -1$ and $b = 1$. For the Unexpected Bernstein bound we take a grid of $\gamma \in \{1/(2b), \dots, 1/(2^k b)\}$ for $k = \lceil \log_2(\sqrt{n/\ln(1/\delta)}/2) \rceil$ and a union bound over the grid, as proposed by Mhammedi et al. (2019). For the split-kl bound we take μ to be the middle value, 0, of the ternary random variable. In the experiments we take $\delta = 0.05$, and truncate the bounds at 1.

In the first experiment, presented in Figure 3.1a, we take $p_{-1} = p_1 = (1 - p_0)/2$ and plot the difference between the values of the bounds and \hat{p} as a function of p_0 . For $p_0 = 0$ the random variable Z is Bernoulli and, as expected, the kl inequality performs the best, followed by split-kl, and then Unexpected Bernstein. As p_0 grows closer to 1, the variance of Z decreases and, also as expected, the kl inequality falls behind, whereas split-kl and Unexpected Bernstein go closely together. Empirical Bernstein falls behind all other bounds throughout most of the range, except slightly outperforming kl when p_0 gets very close to 1.

In the second experiment, presented in Figure 3.1b, we take a skewed random variable with $p_1 = 0.99(1 - p_0)$ and $p_{-1} = 0.01(1 - p_0)$, and again plot the difference between the values of the bounds and \hat{p} as a function of p_0 . This time the kl also starts well for p_0 close to zero, but then falls behind due to its inability of properly handling the values inside the interval. Unexpected Bernstein exhibits the opposite trend due to being based on uncentered second moment, which is high when p_0 is close to zero, even though the variance is small in this case. Empirical Bernstein lags behind all other bounds for most of the range due to poor constants, whereas split-kl matches the tightest bounds, the kl and Unexpected Bernstein, at the endpoints of the range of p_0 , and outperforms all other bounds in the middle of the range, around $p_0 = 0.6$, due to being able to exploit the combinatorics of the problem.

The experiments demonstrate that for ternary random variables the split-kl is a powerful alternative to existing concentration of measure inequalities. To the best of our knowledge, this is also the first empirical evaluation of the Unexpected Bernstein inequality, and it shows that in many cases it is also a powerful inequality. We also observe that in most settings the Empirical Bernstein is weaker than the other three inequalities we consider. Numerical evaluations in additional settings are provided in Appendix 3.6.4.



(a) Comparison of the concentration bounds with $n = 100$, $\delta = 0.05$ and $p_{-1} = p_1 = 0.5(1 - p_0)$.

(b) Comparison of the concentration bounds with $n = 100$, $\delta = 0.05$, $p_1 = 0.99(1 - p_0)$, and $p_{-1} = 0.01(1 - p_0)$.

Figure 3.1: Empirical comparison of the concentration bounds.

3.3 PAC-Bayesian Inequalities

Now we elevate the basic concentration of measure inequalities to the PAC-Bayesian domain. We start with the supervised learning problem setup, then provide a background on existing PAC-Bayesian inequalities, and finish with presentation of the PAC-Bayes-split-kl inequality.

3.3.1 Supervised Learning Problem Setup and Notations

Let \mathcal{X} be a sample space, \mathcal{Y} be a label space, and let $S = \{(X_i, Y_i)\}_{i=1}^n$ be an i.i.d. sample drawn according to an unknown distribution \mathcal{D} on the product-space $\mathcal{X} \times \mathcal{Y}$. Let \mathcal{H} be a hypothesis space containing hypotheses $h : \mathcal{X} \rightarrow \mathcal{Y}$. The quality

of a hypothesis h is measured using the zero-one loss $\ell(h(X), Y) = \mathbf{1}(h(X) \neq Y)$, where $\mathbf{1}(\cdot)$ is the indicator function. The expected loss of h is denoted by $L(h) = \mathbb{E}_{(X,Y) \sim \mathcal{D}}[\ell(h(X), Y)]$, and the empirical loss of h on a sample S is denoted by $\hat{L}(h, S) = \frac{1}{|S|} \sum_{(X,Y) \in S} \ell(h(X), Y)$. We use $\mathbb{E}_{\mathcal{D}}[\cdot]$ as a shorthand for $\mathbb{E}_{(X,Y) \sim \mathcal{D}}[\cdot]$.

PAC-Bayesian bounds bound the generalization error of Gibbs prediction rules. For each input $X \in \mathcal{X}$, Gibbs prediction rule associated with a distribution ρ on \mathcal{H} randomly draws a hypothesis $h \in \mathcal{H}$ according to ρ and predicts $h(X)$. The expected loss of the Gibbs prediction rule is $\mathbb{E}_{h \sim \rho}[L(h)]$ and the empirical loss is $\mathbb{E}_{h \sim \rho}[\hat{L}(h, S)]$. We use $\mathbb{E}_{\rho}[\cdot]$ as a shorthand for $\mathbb{E}_{h \sim \rho}[\cdot]$.

3.3.2 PAC-Bayesian Analysis Background

Now we present a brief background on the relevant results from the PAC-Bayesian analysis.

PAC-Bayes-kl Inequality The PAC-Bayes-kl inequality cited below is one of the tightest known generalization bounds on the expected loss of the Gibbs prediction rule.

Theorem 3.5 (PAC-Bayes-kl Inequality, Seeger, 2002, Maurer, 2004). *For any probability distribution π on \mathcal{H} that is independent of S and any $\delta \in (0, 1)$:*

$$\mathbb{P}\left(\exists \rho \in \mathcal{P} : \text{kl}\left(\mathbb{E}_{\rho}[\hat{L}(h, S)] \parallel \mathbb{E}_{\rho}[L(h)]\right) \geq \frac{\text{KL}(\rho \parallel \pi) + \ln(2\sqrt{n}/\delta)}{n}\right) \leq \delta, \quad (3.6)$$

where \mathcal{P} is the set of all possible probability distributions on \mathcal{H} that can depend on S .

The following relaxation of the PAC-Bayes-kl inequality based on Refined Pinsker's relaxation of the kl divergence helps getting some intuition about the bound (McAllester, 2003). With probability at least $1 - \delta$, for all $\rho \in \mathcal{P}$ we have

$$\mathbb{E}_{\rho}[L(h)] \leq \mathbb{E}_{\rho}[\hat{L}(h, S)] + \sqrt{2\mathbb{E}_{\rho}[\hat{L}(h, S)] \frac{\text{KL}(\rho \parallel \pi) + \ln(2\sqrt{n}/\delta)}{n}} + \frac{2(\text{KL}(\rho \parallel \pi) + \ln(2\sqrt{n}/\delta))}{n}. \quad (3.7)$$

If $\mathbb{E}_{\rho}[\hat{L}(h, S)]$ is close to zero, the middle term in the inequality above vanishes, leading to so-called “fast convergence rates” (convergence of $\mathbb{E}_{\rho}[\hat{L}(h, S)]$ to $\mathbb{E}_{\rho}[L(h)]$ at the rate of $1/n$). However, achieving low $\mathbb{E}_{\rho}[\hat{L}(h, S)]$ is not always possible (Dziugaite and Roy, 2017; Zhou et al., 2019). Subsequent research in PAC-Bayesian analysis has focused on two goals: (1) achieving fast convergence rates when the variance of

prediction errors is low (and not necessarily the errors themselves), and (2) reducing the $\text{KL}(\rho\|\pi)$ term, which may be quite large for large hypothesis spaces. For the first goal Tolstikhin and Seldin (2013) developed the PAC-Bayes-Empirical-Bernstein inequality and Mhammedi et al. (2019) proposed to use excess losses and also derived the alternative PAC-Bayes-Unexpected-Bernstein inequality. For the second goal Ambroladze et al. (2007) suggested to use informed priors and Mhammedi et al. (2019) perfected the idea by proposing to average over “forward” and “backward” construction with informed prior. Next we explain the ideas behind the excess losses and informed priors in more details.

Excess Losses Let h^* be a reference prediction rule that is independent of S . We define the excess loss of a prediction rule h with respect to the reference h^* by

$$\Delta_\ell(h(X), h^*(X), Y) = \ell(h(X), Y) - \ell(h^*(X), Y).$$

If ℓ is the zero-one loss, the excess loss naturally gives rise to ternary random variables, but it is well-defined for any real-valued loss function. We use $\Delta_L(h, h^*) = \mathbb{E}_D[\Delta_\ell(h(X), h^*(X), Y)] = L(h) - L(h^*)$ to denote the expected excess loss of h relative to h^* and $\Delta_{\hat{L}}(h, h^*, S) = \frac{1}{|S|} \sum_{(X, Y) \in S} \Delta_\ell(h(X), h^*(X), Y) = \hat{L}(h) - \hat{L}(h^*)$ to denote the empirical excess loss of h relative to h^* . The expected loss of a Gibbs prediction rule can then be written as

$$\mathbb{E}_\rho[L(h)] = \mathbb{E}_\rho[\Delta_L(h, h^*)] + L(h^*).$$

A bound on $\mathbb{E}_\rho[L(h)]$ can thus be decomposed into a summation of a PAC-Bayes bound on $\mathbb{E}_\rho[\Delta_L(h, h^*)]$ and a bound on $L(h^*)$. When the variance of the excess loss is small, we can use tools that exploit small variance, such as the PAC-Bayes-Empirical-Bernstein, PAC-Bayes-Unexpected-Bernstein, or PAC-Bayes-Split-kl inequalities proposed below, to achieve fast convergence rates for the excess loss. Bounding $L(h^*)$ involves just a single prediction rule and does not depend on the value of $\text{KL}(\rho\|\pi)$. We note that it is essential that the variance and not just the magnitude of the excess loss is small. For example, if the excess losses primarily take values in $\{-1, 1\}$ and average out to zero, fast convergence rates are impossible.

Informed Priors The idea behind informed priors is to split the data into two subsets, $S = S_1 \cup S_2$, and to use S_1 to learn a prior π_{S_1} , and then use it to learn a posterior on S_2 Ambroladze et al. (2007). Note that since the size of S_2 is smaller than the size of S , this approach gains in having potentially smaller $\text{KL}(\rho\|\pi_{S_1})$, but loses in having a smaller sample size in the denominator of the PAC-Bayes bounds.

The balance between the advantage and disadvantage depends on the data: for some data sets it strengthens the bounds, but for some it weakens them. Mhammedi et al. (2019) perfected the approach by proposing to use it in the “forward” and “backward” direction and average over the two. Let S_1 and S_2 be of equal size. The “forward” part uses S_1 to train π_{S_1} and then computes a posterior on S_2 , while the “backward” part uses S_2 to train π_{S_2} and then computes a posterior on S_1 . Finally, the two posteriors are averaged with equal weight and the KL term becomes $\frac{1}{2} (\text{KL}(\rho \parallel \pi_{S_1}) + \text{KL}(\rho \parallel \pi_{S_2}))$. See (Mhammedi et al., 2019) for the derivation.

Excess Losses and Informed Priors Excess losses and informed priors make an ideal combination. If we split S into two equal parts, $S = S_1 \cup S_2$, we can use S_1 to train both a reference prediction rule h_{S_1} and a prior π_{S_1} , and then learn a PAC-Bayes posterior on S_2 , and the other way around. By combining the “forward” and “backward” approaches we can write

$$\mathbb{E}_\rho[L(h)] = \frac{1}{2}\mathbb{E}_\rho[\Delta_L(h, h_{S_1})] + \frac{1}{2}\mathbb{E}_\rho[\Delta_L(h, h_{S_2})] + \frac{1}{2}(L(h_{S_1}) + L(h_{S_2})), \quad (3.8)$$

and we can use PAC-Bayes to bound the first term using the prior π_{S_1} and the data in S_2 , and to bound the second term using the prior π_{S_2} and the data in S_1 , and we can bound $L(h_{S_1})$ and $L(h_{S_2})$ using the “complementary” data in S_2 and S_1 , respectively.

PAC-Bayes-Empirical-Bernstein Inequalities The excess losses are ternary random variables taking values in $\{-1, 0, 1\}$ and, as we have already discussed, the kl inequality is not well-suited for them. PAC-Bayesian inequalities tailored for non-binary random variables were derived by Seldin et al. (2012), Tolstikhin and Seldin (2013), Wu et al. (2021), and Mhammedi et al. (2019). Seldin et al. (2012) derived the PAC-Bayes-Bernstein oracle bound, which assumes knowledge of the variance. Tolstikhin and Seldin (2013) made it into an empirical bound by deriving the PAC-Bayes-Empirical-Bernstein bound for the variance and plugging it into the PAC-Bayes-Bernstein bound of Seldin et al.. Wu et al. (2021) derived an oracle PAC-Bayes-Bennett inequality, which again assumes oracle knowledge of the variance, and showed that it is always at least as tight as the PAC-Bayes-Bernstein, and then also plugged in the PAC-Bayes-Empirical-Bernstein bound on the variance. Mhammedi et al. (2019) derived the PAC-Bayes-Unexpected-Bernstein inequality, which directly uses the empirical second moment. Since we have already shown that the Unexpected Bernstein inequality is tighter than the Empirical Bernstein, and since the approach of Wu et al. requires a combination of two inequalities, PAC-Bayes-Empirical-Bernstein for the variance and PAC-Bayes-Bennett for the

loss, whereas the approach of Mhammedi et al. only makes a single application of PAC-Bayes-Unexpected-Bernstein, we only compare our work to the latter.

We cite the inequality of Mhammedi et al. (2019), which applies to an arbitrary loss function. We use $\tilde{\ell}$ and matching tilde-marked quantities to distinguish it from the zero-one loss ℓ . For any $h \in \mathcal{H}$, let $\tilde{L}(h) = \mathbb{E}_D[\tilde{\ell}(h(X), Y)]$ be the expected tilde-loss of h , and let $\hat{\tilde{L}}(h, S) = \frac{1}{|S|} \sum_{(X,Y) \in S} \tilde{\ell}(h(X), Y)$ be the empirical tilde-loss of h on S .

Theorem 3.6 (PAC-Bayes-Unexpected-Bernstein inequality (Mhammedi et al., 2019)). *Let $\tilde{\ell}(\cdot, \cdot)$ be an arbitrary loss function bounded from above by b for some $b > 0$, and assume that $\hat{\tilde{V}}(h, S) = \frac{1}{|S|} \sum_{(X,Y) \in S} \tilde{\ell}(h(X), Y)^2$ is finite for all h . Let $\psi(u) := u - \ln(1+u)$ for $u > -1$. Then for any distribution π on \mathcal{H} that is independent of S , any $\gamma \in (0, 1/b)$, and any $\delta \in (0, 1)$:*

$$\mathbb{P}\left(\exists \rho \in \mathcal{P} : \mathbb{E}_\rho[\tilde{L}(h)] \geq \mathbb{E}_\rho[\hat{\tilde{L}}(h, S)] + \frac{\psi(-\gamma b)}{\gamma b^2} \mathbb{E}_\rho[\hat{\tilde{V}}(h, S)] + \frac{\text{KL}(\rho||\pi) + \ln \frac{1}{\delta}}{\gamma n}\right) \leq \delta,$$

where \mathcal{P} is the set of all possible probability distributions on \mathcal{H} that can depend on S .

In optimization of the bound, we take a grid of $\gamma \in \{1/(2b), \dots, 1/(2^k b)\}$ for $k = \lceil \log_2(\sqrt{n}/\ln(1/\delta)/2) \rceil$ and a union bound over the grid, as we did for Theorem 3.3.

3.3.3 PAC-Bayes-Split-kl Inequality

Now we present our PAC-Bayes-Split-kl inequality. For an arbitrary loss function $\tilde{\ell}$ taking values in a $[a, b]$ interval for some $a, b \in \mathbb{R}$, let $\tilde{\ell}^+ := \max\{0, \tilde{\ell} - \mu\}$ and $\tilde{\ell}^- := \max\{0, \mu - \tilde{\ell}\}$ for some $\mu \in [a, b]$. For any $h \in \mathcal{H}$, let $\tilde{L}^+(h) = \mathbb{E}_D[\tilde{\ell}^+(h(X), Y)]$ and $\tilde{L}^-(h) = \mathbb{E}_D[\tilde{\ell}^-(h(X), Y)]$. The corresponding empirical losses are denoted by $\hat{\tilde{L}}^+(h, S) = \frac{1}{n} \sum_{i=1}^n \tilde{\ell}^+(h(X_i), Y_i)$ and $\hat{\tilde{L}}^-(h, S) = \frac{1}{n} \sum_{i=1}^n \tilde{\ell}^-(h(X_i), Y_i)$.

Theorem 3.7 (PAC-Bayes-Split-kl Inequality). *Let $\tilde{\ell}(\cdot, \cdot)$ be an arbitrary loss function taking values in a $[a, b]$ interval for some $a, b \in \mathbb{R}$. Then for any distribution π on \mathcal{H} that is independent of S , any $\mu \in [a, b]$, and any $\delta \in (0, 1)$:*

$$\mathbb{P}\left[\exists \rho \in \mathcal{P} : \mathbb{E}_\rho[\tilde{L}(h)] \geq \mu + (b - \mu) \text{kl}^{-1,+} \left(\frac{\mathbb{E}_\rho[\hat{\tilde{L}}^+(h, S)]}{b - \mu}, \frac{\text{KL}(\rho||\pi) + \ln \frac{4\sqrt{n}}{\delta}}{n} \right) - (\mu - a) \text{kl}^{-1,-} \left(\frac{\mathbb{E}_\rho[\hat{\tilde{L}}^-(h, S)]}{\mu - a}, \frac{\text{KL}(\rho||\pi) + \ln \frac{4\sqrt{n}}{\delta}}{n} \right) \right] \leq \delta,$$

where \mathcal{P} is the set of all possible probability distributions on \mathcal{H} that can depend on S .

Proof. We have $\mathbb{E}_\rho[\tilde{L}(h)] = \mu + \mathbb{E}_\rho[\tilde{L}^+(h)] - \mathbb{E}_\rho[\tilde{L}^-(h)]$. Similar to the proof of Theorem 3.4, we take a union bound of PAC-Bayes-kl upper bound on $\mathbb{E}_\rho[\tilde{L}^+(h)]$ and PAC-Bayes-kl lower bound on $\mathbb{E}_\rho[\tilde{L}^-(h)]$. \square

3.3.4 PAC-Bayes-split-kl with Excess Loss and Informed Prior

Looking back at the expected loss decomposition in equation (3.8), we can use PAC-Bayes-split-kl to bound the first two terms and a bound on the binomial tail distribution to bound the last term. For n i.i.d. Bernoulli random variables Z_1, \dots, Z_n with bias $p \in (0, 1)$, we define the binomial tail distribution $\text{Bin}(n, k, p) = \mathbb{P}(\sum_{i=1}^n X_i \leq k)$ and its inverse $\text{Bin}^{-1}(n, k, \delta) = \max\{p : p \in [0, 1] \text{ and } \text{Bin}(n, k, p) \geq \delta\}$. The following theorem relates $\hat{p} = \frac{1}{n} \sum_{i=1}^n Z_i$ and p .

Theorem 3.8 (Test Set Bound (Langford, 2005)). *Let Z_1, \dots, Z_n be n i.i.d. Bernoulli random variables with bias $p \in (0, 1)$ and let $\hat{p} = \frac{1}{n} \sum_{i=1}^n Z_i$ be the empirical mean. Then for any $\delta \in (0, 1)$:*

$$\mathbb{P}(p \geq \text{Bin}^{-1}(n, n\hat{p}, \delta)) \leq \delta.$$

By applying Theorems 3.7 and 3.8 to equation (3.8) we obtain the following result.

Theorem 3.9. *For any $\mu \in [-1, 1]$ and any $\delta \in (0, 1)$:*

$$\mathbb{P}\left(\exists \rho \in \mathcal{P} : \mathbb{E}_\rho[L(h)] \geq \mu + (1 - \mu)(a) - (\mu + 1)(b) + \frac{1}{2}(c)\right) \leq \delta,$$

where \mathcal{P} is the set of all possible probability distributions on \mathcal{H} that can depend on S ,

$$(a) = \text{kl}^{-1,+} \left(\frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_{\hat{L}}^+(h, h_{S_1}, S_2)]}{1 - \mu} + \frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_{\hat{L}}^+(h, h_{S_2}, S_1)]}{1 - \mu}, \frac{\text{KL}(\rho||\pi) + \ln \frac{8\sqrt{n/2}}{\delta}}{n/2} \right),$$

$$(b) = \text{kl}^{-1,-} \left(\frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_{\hat{L}}^-(h, h_{S_1}, S_2)]}{\mu + 1} + \frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_{\hat{L}}^-(h, h_{S_2}, S_1)]}{\mu + 1}, \frac{\text{KL}(\rho||\pi) + \ln \frac{8\sqrt{n/2}}{\delta}}{n/2} \right),$$

in which $\pi = \frac{1}{2}\pi_{S_1} + \frac{1}{2}\pi_{S_2}$, and

$$(c) = \text{Bin}^{-1} \left(\frac{n}{2}, \frac{n}{2} \hat{L}(h_{S_1}, S_2), \frac{\delta}{4} \right) + \text{Bin}^{-1} \left(\frac{n}{2}, \frac{n}{2} \hat{L}(h_{S_2}, S_1), \frac{\delta}{4} \right).$$

The proof is postponed to Appendix 3.6.3.

3.4 Experiments

We evaluate the performance of the PAC-Bayes-split-kl inequality in linear classification and in weighted majority vote using several data sets from UCI and LibSVM repositories (Dua and Graff, 2019; Chang and Lin, 2011). An overview of the data sets is provided in Appendix 3.6.5.1. For linear classification we reproduce the experimental setup of Mhammedi et al. (2019), and for the weighted majority vote we reproduce the experimental setup of Wu et al. (2021).

3.4.1 The Experimental Setup of Mhammedi et al. (2019): Linear Classifiers

In the first experiment we follow the experimental setup of Mhammedi et al. (2019), who consider binary classification problems with linear classifiers in \mathbb{R}^d and Gaussian priors and posteriors. A classifier h_w associated with a vector $w \in \mathbb{R}^d$ makes a prediction on an input X by $h_w(X) = \mathbf{1}(w^\top X > 0)$. The posteriors have the form of Gaussian distributions centered at $w_S \in \mathbb{R}^d$, with covariance Σ_S that depends on a sample S , $\rho = \mathcal{N}(w_S, \Sigma_S)$. The informed priors $\pi_{S_1} = \mathcal{N}(w_{S_1}, \Sigma_{S_1})$ and $\pi_{S_2} = \mathcal{N}(w_{S_2}, \Sigma_{S_2})$ are also taken to be Gaussian distributions centered at w_{S_1} and w_{S_2} , with covariance Σ_{S_1} and Σ_{S_2} , respectively. We take the classifier associated with w_{S_1} as the reference classifier h_{S_1} and the classifier associated with w_{S_2} as the reference classifier h_{S_2} . More details on the construction are provided in Appendix 3.6.5.2.

Figure 3.2 compares the PAC-Bayes-Unexpected-Bernstein bound PBUB and the PAC-Bayes-split-kl bound PBSkl with excess losses and informed priors. The ternary random variables in this setup take values in $\{-1, 0, 1\}$, and we select μ to be the middle value 0. Since the PAC-Bayes-kl bound (PBkl) is one of the tightest known generalization bounds, we take PBkl with informed priors as a baseline. The details on bound calculation and optimization are provided in Appendix 3.6.5.2. In this experiment all the three bounds, PBkl, PBUB, and PBSkl performed comparably. We believe that the reason is that with informed priors the $\text{KL}(\rho||\pi)$ term is small. From the relaxation of the PBkl bound in equation (3.7), we observe that a small $\text{KL}(\rho||\pi)$ term implies smaller difference between fast and slow convergence rates, and thus smaller advantage to bounding the excess loss instead of the raw loss. In other words, we believe that the effect of using informed priors dominates the effect of using excess losses. We note that in order to use excess losses we need to train the reference hypothesis h^* on part of the data and, therefore, training an informed prior on the same data comes at no extra cost.

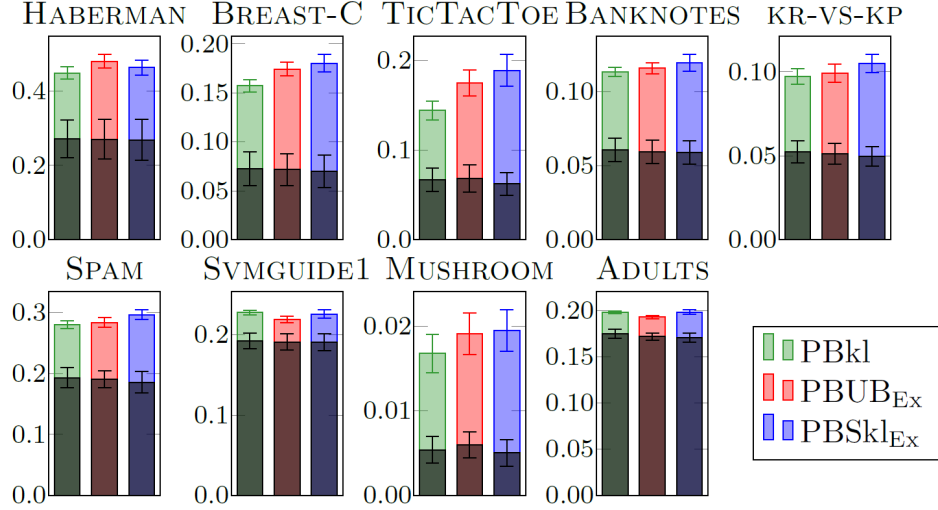


Figure 3.2: Comparison of the bounds and the test losses of the optimized Gaussian posterior ρ^* generated by PBkl with informed priors, PBUB with excess losses and informed priors, and PBSkl with excess losses and informed priors. The test losses of the corresponding bounds are shown in black. We report the mean and the standard deviation over 20 runs of the experiments.

3.4.2 The Experimental Setup of Wu et al. (2021): Weighted Majority Vote

In the second experiment we reproduce the experimental setup of Wu et al. (2021), who consider multiclass classification by a weighted majority vote. Given an input $X \in \mathcal{X}$, a hypothesis space \mathcal{H} , and a distribution ρ on \mathcal{H} , a ρ -weighted majority vote classifier predicts $MV_\rho(X) = \arg \max_{y \in \mathcal{Y}} \mathbb{E}_\rho[\mathbb{1}(h(X) = y)]$. One of the tightest bound for the majority vote is the tandem bound (TND) proposed by Masegosa et al. (2020), which is based on tandem losses for pairs of hypotheses, $\ell(h(X), h'(X), Y) = \mathbb{1}(h(X) \neq Y)\mathbb{1}(h'(X) \neq Y)$, and the second-order Markov's inequality. Wu et al. (2021) proposed two improved forms of the bound, both based on a parametric form of the Chebyshev-Cantelli inequality. The first, CCTND, using Chebyshev-Cantelli with the tandem losses and the PAC-Bayes-kl bound for bounding the tandem losses. The second, CCPBB, using tandem losses with an offset, defined by $\ell_\alpha(h(X), h'(X), Y) = (\mathbb{1}(h(X) \neq Y) - \alpha)(\mathbb{1}(h'(X) \neq Y) - \alpha)$ for $\alpha < 0.5$, and PAC-Bayes-Empirical-Bennett inequality for bounding the tandem losses with an offset. We note that while the tandem losses are binary random variables, tandem losses with an offset are ternary random variables taking values in $\{\alpha^2, -\alpha(1-\alpha), (1-\alpha)^2\}$ and, therefore, application

of Empirical Bernstein type inequalities makes sense. However, in the experiments of Wu et al. CCPBB lagged behind TND and CCTND. We replaced PAC-Bayes-Empirical-Bennett with PAC-Bayes-Unexpected-Bernstein (CCPBUB) and PAC-Bayes-split-kl (CCPBSkl) and showed that the weakness of CCPBB was caused by looseness of PAC-Bayes-Empirical-Bernstein, and that CCPBUB and CCPBSkl lead to tighter bounds that are competitive and sometimes outperforming TND and CCTND. For the PAC-Bayes-split-kl bound we took μ to be the middle value of the tandem loss with an offset, namely, for $\alpha \geq 0$ we took $\mu = \alpha^2$, and for $\alpha < 0$ we took $\mu = -\alpha(1 - \alpha)$.

In Figure 3.3 we present a comparison of the TND, CCTND, CCPBB, CCPBUB, and CCPBSkl bounds on weighted majority vote of heterogeneous classifiers (Linear Discriminant Analysis, k -Nearest Neighbors, Decision Tree, Logistic Regression, and Gaussian Naive Bayes), which adds the two new bounds, CCPBUB and CCPBSkl to the experiment done by Wu et al. (2021). A more detailed description of the experiment and results for additional data sets are provided in Appendix 3.6.5.3. We note that CCPBUB and CCPBSkl consistently outperform CCPBB, demonstrating that they are more appropriate for tandem losses with an offset. The former two bounds perform comparably to TND and CCTND, which operate on tandem losses without an offset. In Appendix 3.6.5.4 we replicate another experiment of Wu et al., where we use the bounds to reweigh trees in a random forest classifier. The results are similar to the results for heterogeneous classifiers.

3.5 Discussion

We have presented the split-kl and PAC-Bayes-split-kl inequalities. The inequalities answer a long-standing open question on how to exploit the structure of ternary random variables in order to provide tight concentration bounds. The proposed split-kl and PAC-Bayes-split-kl inequalities are as tight for ternary random variables, as the kl and PAC-Bayes-kl inequalities are tight for binary random variables.

In our empirical evaluation the split-kl inequality was always competitive with the kl and Unexpected Bernstein inequalities and outperformed both in certain regimes, whereas Empirical Bernstein typically lagged behind. In our experiments in the PAC-Bayesian setting the PAC-Bayes-split-kl was always comparable to PAC-Bayes-Unexpected-Bernstein, whereas PAC-Bayes-Empirical-Bennett most often lagged behind. The first two inequalities were usually comparable to PAC-Bayes-kl, although in some cases the attempt to exploit low variance did not pay off and PAC-Bayes-kl

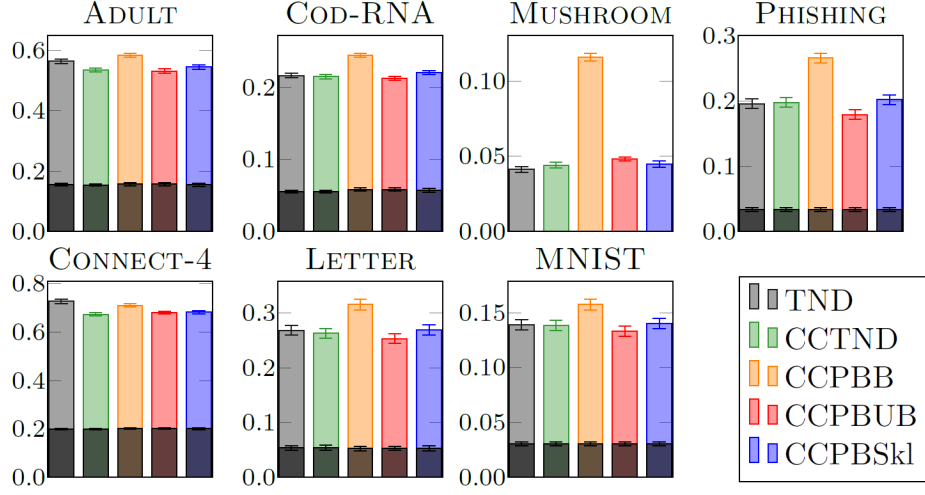


Figure 3.3: Comparison of the bounds and the test losses of the weighted majority vote on ensembles of heterogeneous classifiers with optimized posterior ρ^* generated by TND, CCTND, CCPBB, CCPBUB, and CCPBSkl. The test losses of the corresponding bounds are shown in black. We report the mean and the standard deviation over 10 runs of the experiments.

outperformed, which is also the trend observed earlier by Mhammedi et al. (2019). To the best of our knowledge, this is the first time when the various approaches to exploitation of low variance were directly compared, and the proposed split-kl emerged as a clear winner in the basic setting, whereas in the PAC-Bayes setting in our experiments the PAC-Bayes-Unexpected-Bernstein and PAC-Bayes-split-kl were comparable, and preferable over PAC-Bayes-Empirical-Bernstein and PAC-Bayes-Empirical-Bennett.

3.6 Appendix

3.6.1 Unexpected Bernstein Inequality

3.6.1.1 A Proof of the Unexpected Bernstein Inequality (Theorem 3.3)

The proof is based on the Unexpected Bernstein lemma.

Lemma 3.1 (Unexpected Bernstein lemma (Fan et al., 2015; Mhammedi et al., 2019)). *Let Z_1, \dots, Z_n be i.i.d. random variables bounded from above by $b > 0$, and*

assume that $\sum_{i=1}^n Z_i^2$ is finite. Let $\psi(u) := u - \ln(1 + u)$ for $u \in \mathbb{R}$. Then for any $\gamma \in (0, \frac{1}{b})$:

$$\mathbb{E} \left[e^{\gamma \sum_{i=1}^n (\mathbb{E}[Z_i] - Z_i) - \frac{\psi(-b\gamma)}{b^2} \sum_{i=1}^n Z_i^2} \right] \leq 1.$$

Proof of Theorem 3.3. Recall that by the assumption of the theorem Z_1, \dots, Z_n are i.i.d., bounded from above by $b > 0$, and that $p = \mathbb{E}[Z_i]$ for all i , $\hat{p} = \frac{1}{n} \sum_{i=1}^n Z_i$, and $\hat{\sigma} = \frac{1}{n} \sum_{i=1}^n Z_i^2$. For any $\gamma \in (0, 1/b)$ we have:

$$\begin{aligned} \mathbb{P} \left(p - \hat{p} - \frac{\psi(-\gamma b)}{\gamma b^2} \hat{\sigma} \geq \varepsilon \right) &= \mathbb{P} \left(\sum_{i=1}^n \left(\mathbb{E}[Z_i] - Z_i - \frac{\psi(-\gamma b)}{\gamma b^2} Z_i^2 \right) \geq n\varepsilon \right) \\ &= \mathbb{P} \left(e^{\gamma \sum_{i=1}^n \left(\mathbb{E}[Z_i] - Z_i - \frac{\psi(-\gamma b)}{\gamma b^2} Z_i^2 \right)} \geq e^{\gamma n\varepsilon} \right) \\ &\leq \mathbb{E} \left[e^{\gamma \sum_{i=1}^n \left(\mathbb{E}[Z_i] - Z_i - \frac{\psi(-\gamma b)}{\gamma b^2} Z_i^2 \right)} \right] / e^{\gamma n\varepsilon} \\ &\leq e^{-\gamma n\varepsilon}, \end{aligned}$$

where the first inequality is by application of Markov's inequality and the second inequality is by application of Lemma 3.1. By taking $\delta = e^{-\gamma n\varepsilon}$ and solving for ε we complete the proof. \square

3.6.1.2 A relaxation of the Unexpected Bernstein lemma

We show that a concentration inequality introduced by Cesa-Bianchi et al. (2007) yields a relaxation of the Unexpected Bernstein Lemma. The inequality of Cesa-Bianchi et al. can be used to directly derive a relaxed version of the Unexpected Bernstein lemma, but as we show the result is weaker than the Unexpected Bernstein lemma. Cesa-Bianchi et al. (2007, Lemma 1) have shown that

$$\forall \gamma \geq -1/2 : \quad \gamma - \gamma^2 \leq \ln(1 + \gamma). \quad (3.9)$$

Thus,

$$\forall \gamma \leq 1/2 : \quad -\gamma^2 \leq \gamma + \ln(1 - \gamma) = -\psi(-\gamma). \quad (3.10)$$

This gives a relaxed version of the Unexpected Bernstein lemma. For simplicity, we present it with $b = 1$.

Lemma 3.2 (Relaxed Unexpected Bernstein lemma). *Let Z_1, \dots, Z_n be i.i.d. random variables bounded from above by 1, and assume that $\sum_{i=1}^n Z_i^2$ is finite. Then for any $\gamma \in [0, \frac{1}{2}]$:*

$$\mathbb{E} \left[e^{\gamma \sum_{i=1}^n (\mathbb{E}[Z_i] - Z_i) - \gamma^2 \sum_{i=1}^n Z_i^2} \right] \leq 1.$$

Proof. By (3.10) and Lemma 3.1 we have

$$\mathbb{E} \left[e^{\gamma \sum_{i=1}^n (\mathbb{E}[Z_i] - Z_i) - \gamma^2 \sum_{i=1}^n Z_i^2} \right] \leq \mathbb{E} \left[e^{\gamma \sum_{i=1}^n (\mathbb{E}[Z_i] - Z_i) - \psi(-\gamma) \sum_{i=1}^n Z_i^2} \right] \leq 1.$$

□

We note that it is possible to prove Lemma 3.2 directly by using inequality (3.9) and without using Lemma 3.1, as done by Wintenberger (2017). The first inequality in our proof of Lemma 3.2 shows that it is a relaxation of Lemma 3.1.

3.6.2 A Proof of the PAC-Bayes Unexpected Bernstein Inequality (Theorem 3.6)

The proof is based on using the Unexpected Bernstein lemma within a standard change of measure argument cited in Lemma 2.3 below. We cite the version before the expectations of S' and π are exchanged, which is an intermediate step in the proof of Tolstikhin and Seldin (2013, Lemma 1).

Lemma 3.3 (PAC-Bayes Lemma (Tolstikhin and Seldin, 2013)). *For any function $f_n : \mathcal{H} \times (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathbb{R}$ and for any distribution π on \mathcal{H} , with probability at least $1 - \delta$ over a random draw of S , for all distributions ρ on \mathcal{H} simultaneously:*

$$\mathbb{E}_\rho[f_n(h, S)] \leq \text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta} + \ln \mathbb{E}_{S'}[\mathbb{E}_\pi[e^{f_n(h, S')}]].$$

Proof of Theorem 3.6. Let $f_n(h, S) = \gamma n \left(\tilde{L}(h) - \hat{L}(h, S) \right) - \frac{\psi(-b\gamma)}{b^2} n \hat{\mathbb{V}}(h, S)$. Since π is independent of S by assumption, we can exchange the expectations of S and π . Then by Lemma 3.1 we have $\mathbb{E} [e^{f_n(h, S)}] \leq 1$. By plugging this into Lemma 3.3 and dividing both sides by γn , we complete the proof. □

3.6.3 Proof of Theorem 3.9

To prove the theorem, we need the test set bound (Theorem 3.8), the PAC-Bayes Lemma (Lemma 3.3), and the following lemma.

Lemma 3.4 ((Maurer, 2004)). *Let X_1, \dots, X_n be i.i.d. random variables with mean p and bounded in the $[0, 1]$ interval. Let $\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$ be the empirical mean. Then:*

$$\mathbb{E} [e^{n \text{kl}(\hat{p} \parallel p)}] \leq 2\sqrt{n}.$$

Proof of Theorem 3.9. Recall that

$$\mathbb{E}_\rho[L(h)] = \frac{1}{2}\mathbb{E}_\rho[\Delta_L(h, h_{S_1})] + \frac{1}{2}\mathbb{E}_\rho[\Delta_L(h, h_{S_2})] + \frac{1}{2}(L(h_{S_1}) + L(h_{S_2})). \quad (3.11)$$

First, by applying Theorem 3.8 to $L(h_{S_1})$ and $L(h_{S_2})$, respectively, we have:

$$\mathbb{P}\left(L(h_{S_1}) \geq \text{Bin}^{-1}\left(\frac{n}{2}, \frac{n}{2}\hat{L}(h_{S_1}, S_2), \delta\right)\right) \leq \delta \quad (3.12)$$

and

$$\mathbb{P}\left(L(h_{S_2}) \geq \text{Bin}^{-1}\left(\frac{n}{2}, \frac{n}{2}\hat{L}(h_{S_2}, S_1), \delta\right)\right) \leq \delta. \quad (3.13)$$

Next, since

$$\mathbb{E}_\rho[\Delta_L(h, h_{S_1})] = \mu + \mathbb{E}_\rho[\Delta_L^+(h, h_{S_1})] - \mathbb{E}_\rho[\Delta_L^-(h, h_{S_1})]$$

and

$$\mathbb{E}_\rho[\Delta_L(h, h_{S_2})] = \mu + \mathbb{E}_\rho[\Delta_L^+(h, h_{S_2})] - \mathbb{E}_\rho[\Delta_L^-(h, h_{S_2})],$$

for any $\mu \in [a, b]$ we have

$$\begin{aligned} & \frac{1}{2}\mathbb{E}_\rho[\Delta_L(h, h_{S_1})] + \frac{1}{2}\mathbb{E}_\rho[\Delta_L(h, h_{S_2})] \\ &= \mu + \left(\frac{1}{2}\mathbb{E}_\rho[\Delta_L^+(h, h_{S_1})] + \frac{1}{2}\mathbb{E}_\rho[\Delta_L^+(h, h_{S_2})] \right) - \left(\frac{1}{2}\mathbb{E}_\rho[\Delta_L^-(h, h_{S_1})] + \frac{1}{2}\mathbb{E}_\rho[\Delta_L^-(h, h_{S_2})] \right). \end{aligned} \quad (3.14)$$

Let $\pi = \frac{1}{2}\pi_{S_1} + \frac{1}{2}\pi_{S_2}$, and let S_* be either S_1 or S_2 and $\bar{S}_* = S \setminus S_*$. If h is sampled from π_{S_*} , we take h_{S_*} as a reference hypothesis and estimate the excess loss on \bar{S}_* . Then,

$$\begin{aligned} \mathbb{E}_S \mathbb{E}_\pi \left[e^{\frac{n}{2} \text{kl} \left(\frac{\Delta_L^+(h, h_{S_*}, \bar{S}_*)}{1-\mu} \parallel \frac{\Delta_L^+(h, h_{S_*})}{1-\mu} \right)} \right] &= \frac{1}{2} \sum_{i=1,2} \mathbb{E}_S \mathbb{E}_{\pi_{S_i}} \left[e^{\frac{n}{2} \text{kl} \left(\frac{\Delta_L^+(h, h_{S_i}, \bar{S}_i)}{1-\mu} \parallel \frac{\Delta_L^+(h, h_{S_i})}{1-\mu} \right)} \right] \\ &= \frac{1}{2} \sum_{i=1,2} \mathbb{E}_{S_i} \mathbb{E}_{\pi_{S_i}} \mathbb{E}_{\bar{S}_i} \left[e^{\frac{n}{2} \text{kl} \left(\frac{\Delta_L^+(h, h_{S_i}, \bar{S}_i)}{1-\mu} \parallel \frac{\Delta_L^+(h, h_{S_i})}{1-\mu} \right)} \right] \\ &\leq 2\sqrt{n/2}, \end{aligned}$$

where the second equality is due to the fact that π_{S_*} is independent of \bar{S}_* so they are exchangeable, and the inequality follows by Lemma 3.4.

Therefore, by applying Lemma 3.3 with $f(h, S) = \frac{n}{2} \text{kl} \left(\frac{\Delta_{\hat{L}}^+(h, h_{S_*}, \bar{S}_*)}{1-\mu} \left\| \frac{\Delta_L^+(h, h_{S_*})}{1-\mu} \right\| \right)$, we have with probability at least $1 - \delta$ over S , for all ρ on \mathcal{H} simultaneously:

$$\mathbb{E}_\rho \left[\frac{n}{2} \text{kl} \left(\frac{\Delta_{\hat{L}}^+(h, h_{S_*}, \bar{S}_*)}{1-\mu} \left\| \frac{\Delta_L^+(h, h_{S_*})}{1-\mu} \right\| \right) \right] \leq \text{KL}(\rho \|\pi) + \ln \frac{2\sqrt{n/2}}{\delta}.$$

By the convexity of KL, we further have

$$\text{kl} \left(\mathbb{E}_\rho \left[\frac{\Delta_{\hat{L}}^+(h, h_{S_*}, \bar{S}_*)}{1-\mu} \right] \left\| \mathbb{E}_\rho \left[\frac{\Delta_L^+(h, h_{S_*})}{1-\mu} \right] \right\| \right) \leq \mathbb{E}_\rho \left[\text{kl} \left(\frac{\Delta_{\hat{L}}^+(h, h_{S_*}, \bar{S}_*)}{1-\mu} \left\| \frac{\Delta_L^+(h, h_{S_*})}{1-\mu} \right\| \right) \right],$$

which together gives with probability at least $1 - \delta$ over S , for all ρ on \mathcal{H} simultaneously:

$$\text{kl} \left(\mathbb{E}_\rho \left[\frac{\Delta_{\hat{L}}^+(h, h_{S_*}, \bar{S}_*)}{1-\mu} \right] \left\| \mathbb{E}_\rho \left[\frac{\Delta_L^+(h, h_{S_*})}{1-\mu} \right] \right\| \right) \leq \frac{\text{KL}(\rho \|\pi) + \ln \frac{2\sqrt{n/2}}{\delta}}{n/2}.$$

Similarly, $\Delta_{\hat{L}}^-(h, h_{S_*}, \bar{S}_*)$ and $\Delta_L^-(h, h_{S_*})$ also satisfy with probability at least $1 - \delta$ over S , for all ρ on \mathcal{H} simultaneously:

$$\text{kl} \left(\mathbb{E}_\rho \left[\frac{\Delta_{\hat{L}}^-(h, h_{S_*}, \bar{S}_*)}{\mu+1} \right] \left\| \mathbb{E}_\rho \left[\frac{\Delta_L^-(h, h_{S_*})}{\mu+1} \right] \right\| \right) \leq \frac{\text{KL}(\rho \|\pi) + \ln \frac{2\sqrt{n/2}}{\delta}}{n/2}.$$

Let $\rho = \frac{1}{2}\rho_1 + \frac{1}{2}\rho_2$ be constructed in a similar way to π , where ρ_1 and ρ_2 are probability distributions on \mathcal{H} . If h is sampled from ρ_* , then we take h_{S_*} as a reference hypothesis and estimate the excess loss on \bar{S}_* . In our case, $\rho_1 = \rho_2 = \rho$. Let Δ° denote either Δ^+ or Δ^- . Then,

$$\mathbb{E}_\rho[\Delta_L^\circ(h, h_{S_*})] = \frac{1}{2}\mathbb{E}_\rho[\Delta_L^\circ(h, h_{S_1})] + \frac{1}{2}\mathbb{E}_\rho[\Delta_L^\circ(h, h_{S_2})]$$

and

$$\mathbb{E}_\rho[\Delta_{\hat{L}}^\circ(h, h_{S_*}, \bar{S}_*)] = \frac{1}{2}\mathbb{E}_\rho[\Delta_{\hat{L}}^\circ(h, h_{S_1}, S_2)] + \frac{1}{2}\mathbb{E}_\rho[\Delta_{\hat{L}}^\circ(h, h_{S_2}, S_1)].$$

By taking the inverse of kl, we obtain that with probability at least $1 - \delta$ over S , for all ρ on \mathcal{H} simultaneously:

$$\begin{aligned} & \frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_L^+(h, h_{S_1})]}{1 - \mu} + \frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_L^+(h, h_{S_2})]}{1 - \mu} \\ & \leq \text{kl}^{-1,+} \left(\frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_{\hat{L}}^+(h, h_{S_1}, S_2)]}{1 - \mu} + \frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_{\hat{L}}^+(h, h_{S_2}, S_1)]}{1 - \mu}, \frac{\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{n/2}}{\delta}}{n/2} \right), \end{aligned} \quad (3.15)$$

and with the same probability

$$\begin{aligned} & \frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_L^-(h, h_{S_1})]}{\mu + 1} + \frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_L^-(h, h_{S_2})]}{\mu + 1} \\ & \geq \text{kl}^{-1,-} \left(\frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_{\hat{L}}^-(h, h_{S_1}, S_2)]}{\mu + 1} + \frac{1}{2} \frac{\mathbb{E}_\rho[\Delta_{\hat{L}}^-(h, h_{S_2}, S_1)]}{\mu + 1}, \frac{\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{n/2}}{\delta}}{n/2} \right). \end{aligned} \quad (3.16)$$

Thus, we can bound Eq. (3.14) by Eq.(3.15) and Eq.(3.16). By replacing Eq.(3.11) by the upper bound of each term and taking a union bound, we complete the proof. \square

3.6.4 Empirical Comparison

We present more results on empirical comparison of the concentration inequalities: the kl, the Empirical Bernstein, the Unexpected Bernstein, and the split-kl. In particular, Section 3.6.4.1 expands the empirical comparison in Section 3.2.3 in the body for ternary random variables, and Section 3.6.4.2 studies the empirical comparison of bounded random variables. The source code for replicating the experiments is available at Github².

3.6.4.1 Ternary Random Variables

In this section, we follow the settings and the parameters in Section 3.2.3, considering n i.i.d. samples taking values in $\{-1, 0, 1\}$. For completeness, Figure 3.4b and Figure 3.4c repeats Figures 3.1a and 3.1b while we add Figure 3.4a, where the probability is defined by $p_1 = 0.01(1 - p_0)$ and $p_{-1} = 0.99(1 - p_0)$. In this case, the kl starts well for p_0 close to zero, but similar to the case in Figure 3.4c falls behind due to its inability of properly handling the values inside the interval. The

²<https://github.com/YiShanAngWu/Split-KL-R/tree/main/simulation>

Unexpected Bernstein and the Empirical Bernstein perform similarly when p_0 is small in Figure 3.4c since the bounds are cut to 1, while Unexpected Bernstein falls behind Empirical Bernstein when p_0 is small in Figure 3.4a due to the uncentered second moment. The split-kl matches, and in many cases outperforms, the tightest bounds.

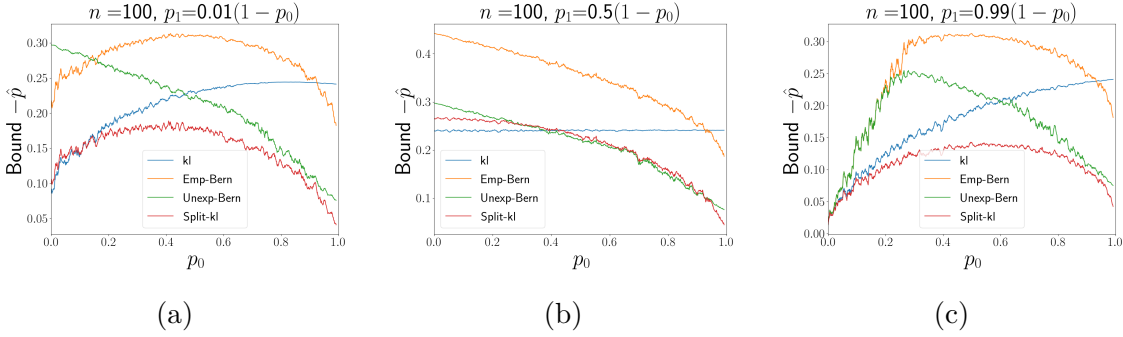


Figure 3.4: Comparison of the concentration bounds with $n = 100$, $\delta = 0.05$, and (a) $p_1 = 0.01(1 - p_0)$ and $p_{-1} = 0.99(1 - p_0)$, (b) $p_{-1} = p_1 = 0.5(1 - p_0)$, (c) $p_1 = 0.99(1 - p_0)$ and $p_{-1} = 0.01(1 - p_0)$.

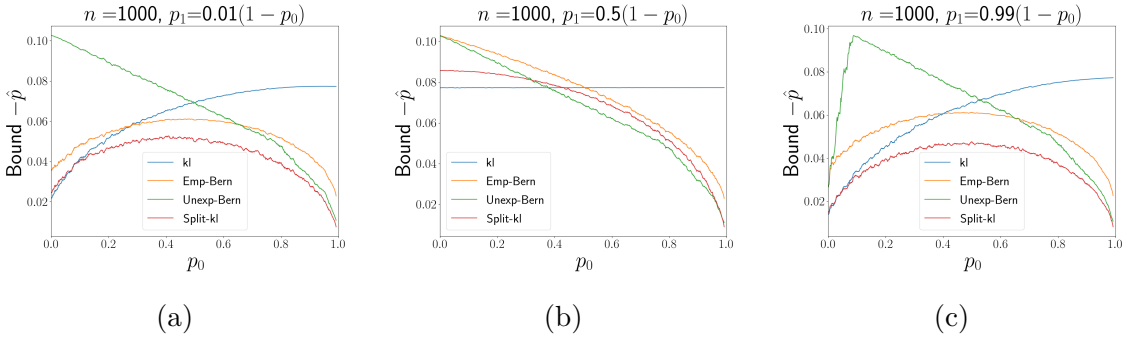


Figure 3.5: Comparison of the concentration bounds with $n = 1000$, $\delta = 0.05$, and (a) $p_1 = 0.01(1 - p_0)$ and $p_{-1} = 0.99(1 - p_0)$, (b) $p_{-1} = p_1 = 0.5(1 - p_0)$, (c) $p_1 = 0.99(1 - p_0)$ and $p_{-1} = 0.01(1 - p_0)$.

Figure 3.5 has the same setting with a larger number of samples $n = 1000$. The trends of the bounds are similar to Figure 3.4. However, the Empirical Bernstein performs better than the Unexpected Bernstein in Figure 3.5a and Figure 3.5c when p_0 is less than 0.6. In both cases, split-kl keeps its leading position. When $p_1 = p_{-1} = (1 - p_0)/2$ (Figure 3.5b), as $p_0 = 0$, the random variable becomes Bernoulli and, as expected,

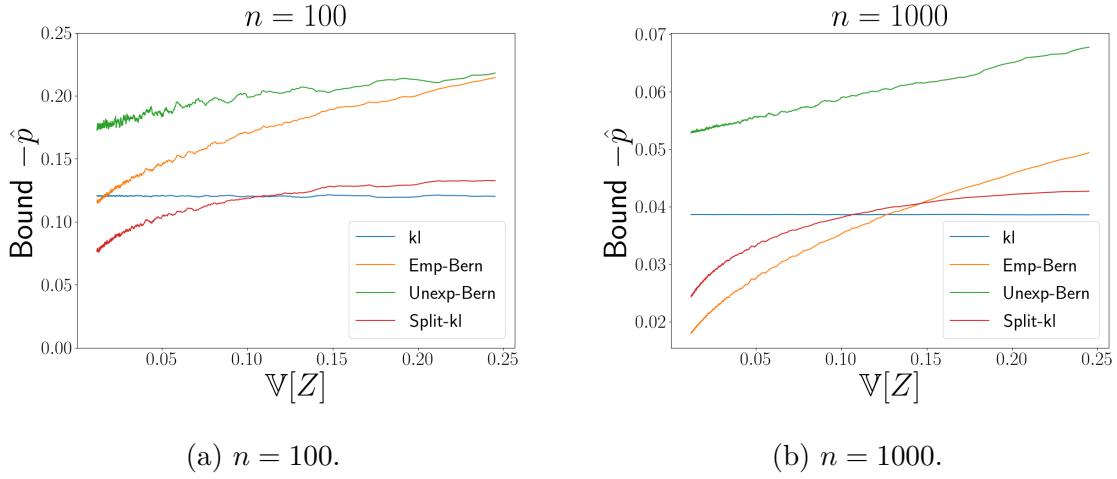


Figure 3.6: Comparison of the concentration bounds for beta distributions with parameters $\alpha = \beta$ taking values in the $[0.01, 10]$ interval, with $\delta = 0.05$, and with the number of samples $n = 100$ and $n = 1000$, respectively.

the kl bound performs the best, followed by the split-kl, and then the two Bernstein bounds. As p_0 grows larger, the kl bound falls behind the other three bounds due to inability of properly handling the values inside the interval. The Unexpected Bernstein, the Empirical Bernstein and the split-kl perform similarly well.

3.6.4.2 Bounded Random Variables

In this section, we study a more general setting, where the i.i.d. random variables Z_1, \dots, Z_n taking values in $[0, 1]$. Naturally, we consider the random variables following beta distribution with parameters $\alpha > 0$ and $\beta > 0$, where the mean $p = \frac{\alpha}{\alpha + \beta}$ and the variance $\mathbb{V}[Z] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$. For the Empirical Bernstein bound, we take $a = 0$ and $b = 1$. For the Unexpected Bernstein bound we take a grid of $\gamma \in \{1/(2b), \dots, 1/(2^k b)\}$ for $b = 1$, $k = \lceil \log_2(\sqrt{n/\ln(1/\delta)}/2) \rceil$, and a union bound over the grid, as in Section 3.2.3. For the split-kl bound we take μ to be the middle value 0.5. Again, in the experiments we take $\delta = 0.05$ and cut the bounds to 1.

In Figure 3.6 we take $\alpha = \beta$ in an interval of $[0.01, 10]$. The mean is a constant $p = 0.5$ throughout the interval and the variance is in an interval of $[0, 0.12, 0.245]$, where a small α and β corresponds to a large variance, and a large α and β corresponds to a small variance. We plot the difference between the values of the bounds and \hat{p} as a function of the variance $\mathbb{V}[Z]$. Since the true mean is a constant, the kl bound

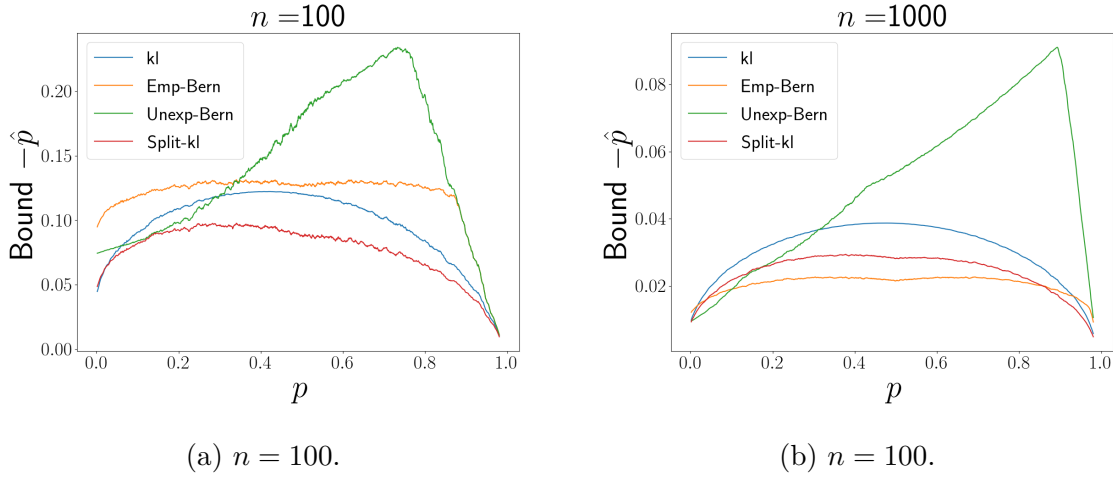


Figure 3.7: Empirical comparison of concentration bounds for beta distribution with parameters α and β and with the number of samples $n = 100$ and $n = 1000$. For $p \in [0, 0.5]$, we take $\beta = 5$ and $\alpha \in [0.01, 5]$ while for $p \in [0.5, 1]$, we take $\alpha = 5$ and $\beta \in [0.01, 5]$.

is also almost a constant throughout the interval. When the variance large, the kl bound performs the best, followed by split-kl and the two Bernstein bounds. When the variance is small, the Empirical Bernstein bound exploit the low variance and outperform all the others when the number of samples is sufficiently large. The Unexpected Bernstein falls behind due the uncentered second moment. The split-kl bound is comparable to the kl bound when the variance is large and also comparable to the tightest bound when the variance is small.

In Figure 3.7 we consider another case where the variances stay similar but the means lie across the spectrum in between 0 and 1. We define the distributions being studied by a combination of two sets of probability distributions. First of all, we take $\beta = 5$ and $\alpha \in [0.01, 5]$, resulting in the mean p in between 0 and 0.5. We define another part by taking $\alpha = 5$ and $\beta \in [0.01, 5]$, resulting in the mean p in between 0.5 and 1. We plot the difference between the values of the bounds and \hat{p} as a function of p . The kl bound is relatively weak around $p = 0.5$ as expected. Since the variances stay similar across the interval, the performance of the Empirical Bernstein stay similar throughout the spectrum, and is tighter than the kl bound when the number of samples is sufficiently large. The split-kl bound is comparable and sometimes outperform the tightest bounds. The Unexpected Bernstein bound again falls behind due to the uncentered second moments.

3.6.5 Experiments

3.6.5.1 Data Sets

As mentioned in Section 3.4, we consider data sets from UCI and LibSVM repositories (Dua and Graff, 2019; Chang and Lin, 2011), as well as Fashion-MNIST from Zalando Research³. An overview of the data sets is listed in Table 3.1, where Banknote stands for Banknote Authentication, Breast-C stands for Breast Cancer Wisconsin, Fashion stands for Fashion-MNIST, Haberman stands for Haberman’s Survival, and kr-vs-kp stands for Chess (King-Rook vs. King-Pawn). For data sets with a training and a testing set, we combine the training and the testing sets.

Linear Classifiers. For the linear classifiers experiments, we consider selected data sets with binary class ($c = 2$). We rescale all the real-valued attributes to the $[-1, 1]$ interval and use one-hot encoding to encode categorical variables to $\{-1, 1\}$, which increases the dimension of the attributes for some of the data sets. In particular, the effective dimension of Adult becomes 108, kr-vs-kp becomes 73, and Mushroom becomes 116. We remove rows containing missing features. For each data set, we shuffle the data sets and take four 5-fold train-test split, which gives 20 runs in total.

Weighted Majority Vote. For the weighted majority vote experiments, including the ensemble of multiple heterogeneous classifiers and the random forest, we consider several binary and multiclass ($c > 2$) data sets. We encode the categorical variables into integers and remove rows containing missing features. For each data set we take 10 runs, and for each run we randomly set aside 20% of sample as the test set.

3.6.5.2 Linear Classifiers

In this section, we describe the details of the experimental setting of the linear classifiers, the details of the bounds, and the details of optimization. The source code for replicating the experiments is available at Github⁴.

Experimental Setting

In this section, we detail the settings and the construction of informed priors and excess losses using linear classifiers with Gaussian posterior. We follow the construction by Mhammedi et al. (2019).

³<https://github.com/zalandoresearch/fashion-mnist>

⁴<https://github.com/YiShanAngWu/Split-KL-R>

Table 3.1: Data set overview. c_{\min} and c_{\max} denote the minimum and maximum class frequency.

Data set	N	d	c	c_{\min}	c_{\max}	Source
Adult	32561	14	2	0.2408	0.7592	LIBSVM (a1a)
Banknote	1372	4	2	0.4447	0.5553	UCI
Breast-C	699	9	2	0.3448	0.6552	UCI
Cod-RNA	59535	8	2	0.3333	0.6667	LIBSVM
Connect-4	67557	126	3	0.0955	0.6583	LIBSVM
Fashion	70000	784	10	0.1000	0.1000	Zalando Research
Haberman	306	3	2	0.2647	0.7353	UCI
kr-vs-kp	3196	36	2	0.48	0.52	UCI
Letter	20000	16	26	0.0367	0.0406	UCI
MNIST	70000	780	10	0.0902	0.1125	LIBSVM
Mushroom	8124	22	2	0.4820	0.5180	LIBSVM
Pendigits	10992	16	10	0.0960	0.1041	LIBSVM
Phishing	11055	68	2	0.4431	0.5569	LIBSVM
Protein	24387	357	3	0.2153	0.4638	LIBSVM
SVMGuide1	3089	4	2	0.3525	0.6475	LIBSVM
SatImage	6435	36	6	0.0973	0.2382	LIBSVM
Sensorless	58509	48	11	0.0909	0.0909	LIBSVM
Shuttle	58000	9	7	0.0002	0.7860	LIBSVM
Spambase	4601	57	2	0.394	0.606	UCI
Splice	3175	60	2	0.4809	0.5191	LIBSVM
TicTacToe	958	9	2	0.347	0.653	UCI
USPS	9298	256	10	0.0761	0.1670	LIBSVM
w1a	49749	300	2	0.0297	0.9703	LIBSVM

As described in Section 3.4, the posterior $\rho = \mathcal{N}(w_S, \Sigma_S)$ is a Gaussian distribution centered at w_S , which is learned on S using regularized logistic regression

$$w_S = \arg \min_{w \in \mathbb{R}^d} \frac{\lambda \|w\|^2}{2} + \frac{1}{|S|} \sum_{(X,Y) \in S} - (Y \ln \phi(w^\top X) + (1 - Y) \ln(1 - \phi(w^\top X))), \quad (3.17)$$

where $\phi(x) := 1/(1 + e^{-x})$ for $x \in \mathbb{R}$ is the sigmoid function. The covariance of the posterior is a diagonal matrix $\sigma^2 I_d$, where the variance σ^2 is learned from the corresponding PAC-Bayes bounds. We use the informed priors in all the PAC-Bayes

bounds. The informed priors $\pi_{S_1} = \mathcal{N}(w_{S_1}, \Sigma_{S_1})$ and $\pi_{S_2} = \mathcal{N}(w_{S_2}, \Sigma_{S_2})$ are also chosen to be Gaussian distributions over \mathbb{R}^d , where the centers of the distributions are learned similarly using regularized logistic regression on the corresponding sample S_1 and S_2 . If using excess losses, we take the classifier associated with w_{S_1} as the reference classifier h_{S_1} for the “forward” approach and take the classifier associated with w_{S_2} as the reference classifier h_{S_2} for the “backward” approach. We let the covariance of the informed priors to be also diagonal matrices $\Sigma_{S_1} = \Sigma_{S_2} = \sigma_\pi^2 I_d$, where σ_π^2 is selected from a grid $\mathcal{G} = \{1/2, \dots, 1/2^j\}$ for $j = \lceil \log_2 |S| \rceil$.

For all data sets, we use $\lambda = 0.01$ in equation (3.17) and solve it using the BFGS algorithm. For all the bounds, we take $\delta = 0.05$. Note that to be able to select the variance of the priors from a grid \mathcal{G} , we have to take a union bound over \mathcal{G} . Since the hypothesis space is infinitely large, we approximate the excess risk by drawing 100 classifiers from the posterior ρ and compute the excess losses with respect to the reference classifiers.

Bounds

As mentioned in the body that we used informed priors for all the bounds we applied. The PBSkl_{Ex} bound is presented in Theorem 3.9, while the PBUB_{Ex} bound and the PBkl bound will be presented in the following. The idea to derive PAC-Bayes bounds with informed priors in general is similar to the technique used in the proof of Theorem 3.9 in Appendix 3.6.3.

The key element of the derivations is to bound $\mathbb{E}_{S'}[\mathbb{E}_\pi[e^{f_n(h, S')}]]$ for a given function $f_n : \mathcal{H} \times (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathbb{R}$ in Lemma 2.3. Let the prior $\pi = \frac{1}{2}\pi_{S_1} + \frac{1}{2}\pi_{S_2}$, and let S_* be either S_1 or S_2 . If h is sampled from π_{S_*} , we estimate the loss on $\bar{S}_* = S \setminus S_*$. Then,

$$\mathbb{E}_S \mathbb{E}_\pi [e^{f_n(h, S)}] = \frac{1}{2} \sum_{i=1,2} \mathbb{E}_S \mathbb{E}_{\pi_{S_i}} [e^{f_n(h, S)}] = \frac{1}{2} \sum_{i=1,2} \mathbb{E}_{S_i} \mathbb{E}_{\pi_{S_i}} \mathbb{E}_{\bar{S}_i} [e^{f_n(h, S)}],$$

where the second equality is due to the fact that π_{S_*} is independent of \bar{S}_* so they are exchangeable. We will then select the function $f_n(h, S)$ later such that $\mathbb{E}_{S_i} \mathbb{E}_{\pi_{S_i}} \mathbb{E}_{\bar{S}_i} [e^{f_n(h, S)}]$ is bounded for $i = 1, 2$.

Similarly, we let $\rho = \frac{1}{2}\rho_1 + \frac{1}{2}\rho_2$. If h is sampled from ρ_* , we estimate the loss on $\bar{S}_* = S \setminus S_*$. Then we have

$$\mathbb{E}_\rho[\tilde{L}(h)] = \frac{1}{2}\mathbb{E}_{\rho_1}[\tilde{L}(h)] + \frac{1}{2}\mathbb{E}_{\rho_2}[\tilde{L}(h)] \quad (3.18)$$

and

$$\mathbb{E}_\rho[\hat{L}(h, S_*)] = \frac{1}{2}\mathbb{E}_{\rho_1}[\hat{L}(h, S_2)] + \frac{1}{2}\mathbb{E}_{\rho_2}[\hat{L}(h, S_1)] \quad (3.19)$$

for any loss $\tilde{\ell}$ and the corresponding quantities following the definitions in Section 3.3.3. We assume that $\rho_1 = \rho_2 = \rho$ in all the bounds. Note that for simpler computation, we replace $\text{kl}(\rho\|\pi)$ by its upper bound $\frac{1}{2}\text{kl}(\rho\|\pi_{S_1}) + \frac{1}{2}\text{kl}(\rho\|\pi_{S_2})$ for all the bounds in the experiments.

PAC-Bayes-kl bound with Informed Priors (PBkl). We take the PAC-Bayes-kl bound with informed priors as the baseline:

$$\mathbb{E}_\rho[L(h)] \leq \text{kl}^{-1,+} \left(\frac{1}{2}\mathbb{E}_\rho[\hat{L}(h, S_1)] + \frac{1}{2}\mathbb{E}_\rho[\hat{L}(h, S_2)], \frac{\text{KL}(\rho\|\pi) + \ln \frac{2|\mathcal{G}|\sqrt{n/2}}{\delta}}{n/2} \right),$$

which is obtained by letting $f_n(h, S) = \frac{n}{2}\text{kl}(\hat{L}(h, \bar{S}_*)\|L(h))$ and plugging it into Lemma 2.3. In particular, we have $\mathbb{E}_{S_i}\mathbb{E}_{\pi_{S_i}}\mathbb{E}_{\bar{S}_i}[e^{f_n(h, S)}] = \mathbb{E}_{S_i}\mathbb{E}_{\pi_{S_i}}\mathbb{E}_{\bar{S}_i}[e^{\frac{n}{2}\text{kl}(\hat{L}(h, \bar{S}_i)\|L(h))}] \leq 2\sqrt{n/2}$ for $i = 1, 2$ by Lemma 3.4. Also, by the convexity of KL, we further have

$$\text{kl} \left(\mathbb{E}_\rho[\hat{L}(h, S_*)] \|\mathbb{E}_\rho[L(h)] \right) \leq \mathbb{E}_\rho \left[\text{kl}(\hat{L}(h, S_*)\|L(h)) \right].$$

By taking the inverse of kl , applying the relations in Eq. (3.18) and Eq. (3.19), and taking a union bound over \mathcal{G} , we obtain the desired formula.

PAC-Bayes-Unexpected Bernstein Bound with Excess Loss and Informed Priors (PBUB_{Ex}). Let $\Delta_{\hat{\psi}}(h, h^*, S) = \frac{1}{|S|} \sum_{(X, Y) \in S} (\Delta_\ell(h(X), h^*(X), Y))^2$ denote the average of the second moment of the excess losses. Then, the PBUB_{Ex} has the form:

$$\begin{aligned} \mathbb{E}_\rho[L(h)] &\leq \frac{1}{2}\mathbb{E}_\rho[\Delta_{\hat{L}}(h, h_{S_1}, S_2)] + \frac{1}{2}\mathbb{E}_\rho[\Delta_{\hat{L}}(h, h_{S_2}, S_1)] \\ &+ \frac{\psi(-\gamma b)}{\gamma b^2} \left(\frac{1}{2}\mathbb{E}_\rho[\Delta_{\hat{\psi}}(h, h_{S_1}, S_2)] + \frac{1}{2}\mathbb{E}_\rho[\Delta_{\hat{\psi}}(h, h_{S_2}, S_1)] \right) + \frac{\text{KL}(\rho\|\pi) + \ln \frac{3|\mathcal{G}||\Gamma|}{\delta}}{\gamma(n/2)} \\ &+ \text{Bin}^{-1} \left(\frac{n}{2}, \frac{n}{2}\hat{L}(h_{S_1}, S_2), \frac{\delta}{3|\mathcal{G}|} \right) + \text{Bin}^{-1} \left(\frac{n}{2}, \frac{n}{2}\hat{L}(h_{S_2}, S_1), \frac{\delta}{3|\mathcal{G}|} \right), \end{aligned}$$

where $|\Gamma|$ comes from a union bound over a grid of $\gamma \in \Gamma = \{1/(2b), \dots, 1/(2^k b)\}$ for $k = \lceil \log_2(\sqrt{|S|}/\ln(1/\delta))/2 \rceil$ when applying the PAC-Bayes-Unexpected-Bernstein inequality.

The last line of the bound is by applying Theorem 3.8 to $L(h_{S_1})$ and $L(h_{S_2})$ as in Theorem 3.9, while the first two lines of the bound are derived from applying the PAC-Bayes-Unexpected-Bernstein inequality to the first two terms in equation (3.8). In particular, let $f_n(h, S) = \gamma \frac{n}{2} (\Delta_L(h, h_{S_*}) - \Delta_{\hat{L}}(h, h_{S_*}, \bar{S}_*)) - \frac{\psi(-b\gamma)}{b^2} \frac{n}{2} \Delta_{\hat{\Psi}}(h, h_{S_*}, \bar{S}_*)$ and plug it into Lemma 2.3. Then we have

$$\mathbb{E}_{S_i} \mathbb{E}_{\pi_{S_i}} \mathbb{E}_{\bar{S}_i} [e^{f_n(h, S)}] = \mathbb{E}_{S_i} \mathbb{E}_{\pi_{S_i}} \mathbb{E}_{\bar{S}_i} [e^{\gamma \frac{n}{2} (\Delta_L(h, h_{S_i}) - \Delta_{\hat{L}}(h, h_{S_i}, \bar{S}_i)) - \frac{\psi(-b\gamma)}{b^2} \frac{n}{2} \Delta_{\hat{\Psi}}(h, h_{S_i}, \bar{S}_i)}] \leq 1$$

for $i = 1, 2$ by Lemma 3.1. By moving the empirical quantities to the right hand side, applying the relations in Eq. (3.18) and Eq. (3.19), and taking the union bounds, we obtain the desired formula.

PAC-Bayes-spli-kl Bound with Excess Loss and Informed Priors (PBSkl_{Ex}).

The bound is stated in Theorem 3.9, except that we replace δ by $\delta/|\mathcal{G}|$ for the union bound of \mathcal{G} . We take $\mu = 0$ for the bound in the experiments.

Optimization

Since the center of the posterior w_S is learned using regularized logistic regression, the only thing remains is to decide the variance of the posterior σ^2 using the PAC-Bayes bounds. In general, the variance can be any non-negative values since the bound holds with high probability for all ρ simultaneously. For simpler computation, we only consider the variance taking the same value as the variance of the priors *i.e.*, taking $\sigma^2 = \sigma_\pi^2 \in \mathcal{G}$. For each PAC-Bayes bounds, we find the optimal σ^2 by iterating over variances $\sigma^2 = \sigma_\pi^2 \in \mathcal{G}$ and return the one that corresponds to the tightest bound. We approximate $\mathbb{E}_\rho[\cdot]$ by sampling 100 classifiers from ρ .

The inverse kl in the PBkl and the PBSkl_{Ex} bounds can be computed by binary search. The inverse of the binomial tail distribution in the PBUB_{Ex} and the PBSkl_{Ex} bounds can also be computed by binary search. To optimize the PBUB_{Ex} bound, we also need to iterate over $\gamma \in \Gamma$.

3.6.5.3 Ensemble of Multiple Heterogeneous Classifiers

In this section, we describe the details of the experimental setting of the ensemble of multiple heterogeneous classifiers, the details of bounds and optimization, and lastly, the results. The source code for replicating the experiments is available at Github⁵.

⁵<https://github.com/StephanLorenzen/MajorityVoteBounds>

Experimental Setting

In this experiment, we follow the setting in Wu et al. (2021). We take the following standard classifiers available in *scikit-learn* using default parameters to build the ensemble: 1. **Linear Discriminant Analysis** 2. **Decision Tree** 3. **Logistic Regression** 4. **Gaussian Naive Bayes**. We also take three versions of **k-Nearest Neighbors**: 1. $k = 3$ with uniform weights (*i.e.*, all points in each neighborhood are weighted equally) 2. $k = 5$ with uniform weights, and 3. $k = 5$ with the weights of the points are defined by the inverse of their L2 distance. Thus, there are 7 classifiers for ensemble in total.

Ensemble Construction by Bagging. We follow the construction used by Masegosa et al. (2020); Wu et al. (2021). For each classifier h , we generate a random split of the data set S into a pair of subsets $S = S_h \cup \bar{S}_h$, where $\bar{S}_h = S \setminus S_h$. We generate the split by the standard bagging method, where S_h contains $0.8|S|$ samples randomly subsampled with replacement from S . We train the classifier on S_h , and estimate the expected loss on the out-of-bag (OOB) sample \bar{S}_h to make an unbiased estimation. The resulting set of classifiers produces an ensemble, while the estimates are used for calculating the bounds and deciding the weights of the ensemble. In particular, we estimate the expected loss by $\hat{L}(h, \bar{S}_h)$, and let $n = \min_h |\bar{S}_h|$. In the remaining of the paper, we call the tandem loss with an offset α by the α -tandem loss. Then for a pair of classifiers h and h' , we take the overlap of the OOB sample $\bar{S}_h \cap \bar{S}_{h'}$ to estimate the unbiased tandem loss $\hat{L}(h, h', \bar{S}_h \cap \bar{S}_{h'})$, α -tandem loss $\hat{L}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})$, the second moment of the α -tandem loss $\hat{V}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})$, the variance of the α -tandem loss $\hat{\text{Var}}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})$, as well as the splits of the α -tandem loss $\hat{L}_\alpha^+(h, h', \bar{S}_h \cap \bar{S}_{h'})$ and $\hat{L}_\alpha^-(h, h', \bar{S}_h \cap \bar{S}_{h'})$. Let $m = \min_{h, h'} |\bar{S}_h \cap \bar{S}_{h'}|$ be the minimum size of the overlap.

Bounds and Optimization

The bounds we are comparing in this section are the TND, CCTND, CCPBB, CCPBUB, and CCPBSkl bounds. The derivations of the first three bounds are provided in Masegosa et al. (2020); Wu et al. (2021), while we will provide the derivations of the CCPBUB bound and the CCPBSkl bound.

In the experiments, we take $\delta = 0.05$ and take π to be a uniform distribution over the classifiers. For CCPBB, CCPBUB and CCPBSkl bounds, we take a grid of $\alpha \in [-0.5, 0.5]$ since the bounds are not differentiable w.r.t α . Note that we don't need a union bound over α (Wu et al., 2021). To optimize the weighting ρ , we applied iRProp+ for the gradient based optimization (Igel and Hüsken, 2003; Florescu and

Igel, 2018), until the bound did not improve more than 10 for 10 iterations. To find the optimal ρ and the parameters, we start by $\rho = \pi$, and apply alternating minimization until the bound doesn't change for more than 10^{-9} . The details of alternating minimization for each bound are provided below.

We first cite the three existing bounds.

Tandem Bound (TND) (Masegosa et al., 2020) They used the following formula to compute the bound after obtaining the optimal weights ρ :

$$L(\text{MV}_\rho) \leq 4 \text{kl}^{-1,+} \left(\mathbb{E}_{\rho^2}[\hat{L}(h, h', \bar{S}_h \cap \bar{S}_{h'})], \frac{2 \text{KL}(\rho \|\pi) + \ln(4\sqrt{m}/\delta)}{m} \right),$$

and used the following relaxation, based on the PAC-Bayes- λ inequality 3.6.6, for easier optimization:

$$L(\text{MV}_\rho) \leq 4 \left(\frac{\mathbb{E}_{\rho^2}[\hat{L}(h, h', \bar{S}_h \cap \bar{S}_{h'})]}{1 - \lambda/2} + \frac{2 \text{KL}(\rho \|\pi) + \ln(2\sqrt{m}/\delta)}{\lambda(1 - \lambda/2)m} \right) \quad (3.20)$$

for any $\lambda \in (0, 2)$. The bound can be optimized by implementing alternating minimization: Given ρ , starting with $\rho = \pi$, find the corresponding optimal λ (Sec. 3.6.6). Then given λ , optimize ρ by projected gradient descent.

Chebyshev-Cantelli bound with TND empirical loss estimate bound (CCTND) (Wu et al., 2021) They used the following formula to compute the bound after obtaining the optimal weights ρ :

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \alpha)^2} \left[\text{kl}^{-1,+} \left(\mathbb{E}_{\rho^2}[\hat{L}(h, h', \bar{S}_h \cap \bar{S}_{h'})], \frac{2 \text{KL}(\rho \|\pi) + \ln(4\sqrt{m}/\delta)}{m} \right) - 2\alpha \text{kl}^{-1,\circ} \left(\mathbb{E}_\rho[\hat{L}(h, \bar{S}_h)], \frac{\text{KL}(\rho \|\pi) + \ln(4\sqrt{n}/\delta)}{n} \right) + \alpha^2 \right],$$

for $\alpha < 0.5$, where \circ is “−” for $\alpha \geq 0$ and “+” otherwise. On the other hand, they used the following relaxations, based on the PAC-Bayes- λ inequality 3.6.6, for easier optimization:

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \alpha)^2} \left[\frac{\mathbb{E}_{\rho^2}[\hat{L}(h, h', \bar{S}_h \cap \bar{S}_{h'})]}{1 - \frac{\lambda}{2}} + \frac{2 \text{KL}(\rho \|\pi) + \ln(4\sqrt{m}/\delta)}{\lambda(1 - \frac{\lambda}{2})m} - 2\alpha \left(\left(1 - \frac{\gamma}{2}\right) \mathbb{E}_\rho[\hat{L}(h, \bar{S}_h)] - \frac{\text{KL}(\rho \|\pi) + \ln(4\sqrt{n}/\delta)}{\gamma n} \right) + \alpha^2 \right]$$

for $0 \leq \alpha < 0.5$, and

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \alpha)^2} \left[\frac{\mathbb{E}_{\rho^2}[\hat{L}(h, h', \bar{S}_h \cap \bar{S}_{h'})]}{1 - \frac{\lambda}{2}} + \frac{2 \text{KL}(\rho \parallel \pi) + \ln(4\sqrt{m}/\delta)}{\lambda \left(1 - \frac{\lambda}{2}\right) m} \right. \\ \left. - 2\alpha \left(\frac{\mathbb{E}_{\rho}[\hat{L}(h, \bar{S}_h)]}{1 - \frac{\gamma}{2}} + \frac{\text{KL}(\rho \parallel \pi) + \ln(4\sqrt{n}/\delta)}{\gamma \left(1 - \frac{\gamma}{2}\right) n} \right) + \alpha^2 \right]$$

for $\alpha < 0$. The optimization of the bound can, again, be done by alternating minimization of the following steps: 1. Given α and ρ , where we start with $\alpha = 0$ and $\rho = \pi$, compute the corresponding closed-form minimizer λ and γ (Sec. 3.6.6). 2. Given ρ , λ and γ , find the closed-form minimizer α . 3. Given parameters α , λ , and γ , optimize over ρ using projected gradient descent.

Chebyshev-Cantelli bound with PAC-Bayes-Bennett loss estimate bound (CCPBB) (Wu et al., 2021) The CCPBB bound has the following formula for both computing the bound and optimization:

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \alpha)^2} \left[\mathbb{E}_{\rho^2}[\hat{L}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})] + \frac{2 \text{KL}(\rho \parallel \pi) + \ln \frac{2k_\lambda k_\gamma}{\delta}}{\gamma m} \right. \\ \left. + \frac{\phi(\gamma K_\alpha)}{\gamma K_\alpha^2} \left(\frac{\mathbb{E}_{\rho^2}[\hat{\text{Var}}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})]}{1 - \frac{\lambda m}{2(m-1)}} + \frac{K_\alpha^2 \left(2 \text{KL}(\rho \parallel \pi) + \ln \frac{2k_\lambda k_\gamma}{\delta}\right)}{n \lambda \left(1 - \frac{\lambda m}{2(m-1)}\right)} \right) \right],$$

where $\phi(x) = e^x - x - 1$ and $K_\alpha = \max\{1 - \alpha, 1 - 2\alpha\}$ is the length of the range of the α -tandem loss. The parameter γ is taken in a grid $\{\gamma_1, \dots, \gamma_{k_\gamma}\}$, where $\gamma_i > 0$ for all i and λ is taken in a grid $\{\lambda_1, \dots, \lambda_{k_\lambda}\}$, where $\lambda_i \in \left(0, \frac{2(n-1)}{n}\right)$ for all i . k_γ and k_λ in the bound come from the union bounds over a grid of γ and a grid of λ .

To optimize the bound, we take a grid of $\alpha \in [-0.5, 0.5]$ and iterate over α in the grid. For a given α , we first compute $\hat{L}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})$ and $\hat{\text{Var}}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})$ for all h, h' . Then, optimize the bound for a fix α by alternating the following steps: 1. Given ρ , starting with $\rho = \pi$, find the corresponding optimal λ , and then the optimal γ in the grids. 2. Given λ and γ , optimize ρ by projected gradient descent.

Next, we present the two new bounds CCPBUB and CCPBSkl, which are based on the oracle parametric form of the Chebyshev-Cantelli bound (Wu et al., 2021, Theorem 8): For all ρ and for all $\alpha < 0.5$

$$L(\text{MV}_\rho) \leq \frac{\mathbb{E}_{\rho^2}[L_\alpha(h, h')]}{(0.5 - \alpha)^2}. \quad (3.21)$$

By applying the PAC-Bayes-Unexpected-Bernstein inequality to the α -tandem loss, we obtain the CCPBUB bound, while by applying the PAC-Bayes-split-kl inequality to the α -tandem loss, we obtain the CCPBSkl bound.

Chebyshev-Cantelli bound with PAC-Bayes-Unexpected-Bernstein loss estimate bound (CCPBUB) By applying the PAC-Bayes-Unexpected-Bernstein inequality to the α -tandem loss in equation (3.21), with the upper bound of the α -tandem loss $b = (1 - \alpha)^2$ for $\alpha < 0.5$, we obtain the bound:

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \alpha)^2} \left[\mathbb{E}_{\rho^2}[\hat{L}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})] + \frac{\psi(-\gamma(1 - \alpha)^2)}{\gamma(1 - \alpha)^4} \mathbb{E}_{\rho^2}[\hat{V}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})] + \frac{2 \text{KL}(\rho \parallel \pi) + \ln \frac{k_\gamma}{\delta}}{\gamma m} \right],$$

where the 2 in front of KL comes from the fact that $\text{KL}(\rho^2 \parallel \pi^2) = 2 \text{KL}(\rho \parallel \pi)$. As in the previous experiments, we take a grid of $\gamma \in \{1/(2(1 - \alpha)^2), \dots, 1/(2^{k_\gamma}(1 - \alpha)^2)\}$ for $k_\gamma = \lceil \log_2(\sqrt{m}/\ln(1/\delta))/2 \rceil$ when applying the PAC-Bayes-Unexpected-Bernstein inequality.

Similar to the optimization of the CCPBB bound, we again take a grid of $\alpha \in [-0.5, 0.5]$ and iterate over α in the grid. For a given α , we first compute $\hat{L}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})$ and $\hat{V}_\alpha(h, h', \bar{S}_h \cap \bar{S}_{h'})$ for all h, h' . Then, optimize the bound for a fix α by alternating minimization of ρ and γ in the grid. We initialize $\rho = \pi$ and optimize it by projected gradient descent.

Chebyshev-Cantelli bound with PAC-Bayes-split-kl loss estimate bound (CCPBSkl) Similarly, by applying the PAC-Bayes-split-kl inequality to the α -tandem loss in equation (3.21), we obtain the following formula to compute the bound after obtaining the optimal weights ρ :

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \alpha)^2} \left[\mu + (b - \mu) \text{kl}^{-1,+} \left(\frac{\mathbb{E}_{\rho^2}[\hat{L}_\alpha^+(h, h', \bar{S}_h \cap \bar{S}_{h'})]}{b - \mu}, \frac{2 \text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta}}{m} \right) - (\mu - a) \text{kl}^{-1,-} \left(\frac{\mathbb{E}_{\rho^2}[\hat{L}_\alpha^-(h, h', \bar{S}_h \cap \bar{S}_{h'})]}{\mu - a}, \frac{2 \text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta}}{m} \right) \right].$$

We use the following relaxation formula, which is based on the PAC-Bayes- λ inequality 3.6.6, for optimization.

$$L(\text{MV}_\rho) \leq \frac{1}{(0.5 - \alpha)^2} \left[\mu + (b - \mu) \left(\frac{\mathbb{E}_{\rho^2}[\hat{L}_\alpha^+(h, h', \bar{S}_h \cap \bar{S}_{h'})]}{(b - \mu)(1 - \lambda/2)} + \frac{2 \text{KL}(\rho \|\pi) + \ln \frac{4\sqrt{m}}{\delta}}{\lambda(1 - \lambda/2)m} \right) - (\mu - a) \left(\left(1 - \frac{\gamma}{2}\right) \frac{\mathbb{E}_{\rho^2}[\hat{L}_\alpha^-(h, h', \bar{S}_h \cap \bar{S}_{h'})]}{\mu - a} - \frac{2 \text{KL}(\rho \|\pi) + \ln \frac{4\sqrt{m}}{\delta}}{\gamma m} \right) \right].$$

$2 \text{KL}(\rho \|\pi)$ in both bounds again comes from $\text{KL}(\rho^2 \|\pi^2) = 2 \text{KL}(\rho \|\pi)$. Recall that the α -tandem loss takes values in $\{(1 - \alpha)^2, -\alpha(1 - \alpha), \alpha^2\}$. For $\alpha < 0.5$, $(1 - \alpha)^2$ has the largest value. Therefore, we take $b = (1 - \alpha)^2$ in the bound. Furthermore, for $\alpha < 0$ we have $\alpha^2 < -\alpha(1 - \alpha)$, and for $\alpha \geq 0$ we have $-\alpha(1 - \alpha) \leq \alpha^2$. Therefore, for $\alpha < 0$, we take $a = \alpha^2$ and μ to be the middle value $-\alpha(1 - \alpha)$, while for $\alpha \geq 0$, we take $a = -\alpha(1 - \alpha)$ and $\mu = \alpha^2$.

To optimize the bound, we again take a grid of $\alpha \in [-0.5, 0.5]$ and iterate over α in the grid. For a given α , we compute the parameters a, b, μ , and then compute the losses $\hat{L}_\alpha^+(h, h', \bar{S}_h \cap \bar{S}_{h'})$ and $\hat{L}_\alpha^-(h, h', \bar{S}_h \cap \bar{S}_{h'})$ for all h, h' . The optimization of the bound for a fixed α can be done by alternating minimization: 1. Given ρ , starting with $\rho = \pi$, compute the corresponding optimal λ and γ (Sec. 3.6.6). 2. Given λ and γ , optimize over ρ using projected gradient descent.

Results

We presented in the body the results of the selected data sets. Here we show the results on more data sets. We present the results for binary data sets in Figure 3.8, while we present the results for multiclass data sets in Figure 3.9. Taking $\alpha = 0$ for CCTND, CCPBB, CCPBUB, and CCPBSkl bounds collapses to the TND bound. Therefore, we take TND as a baseline. In both figures, CCTND performs similar to, and often better than the baseline. The second bound, CCPBB, using the α -tandem loss, lags behind due to nested application of concentration bounds. The two new bounds based the α -tandem loss, CCPBUB and CCPBSkl, clearly improve the shortage of the CCPBB bound and often provide tighter bounds than the rest.

3.6.5.4 Random Forest Majority Vote Classifiers

In this section, we describe the details of the experimental setting of random forest and the results of the experiments. Since both the ensemble of the heterogeneous classifiers and the random forest are examples of weighted majority vote, the bounds

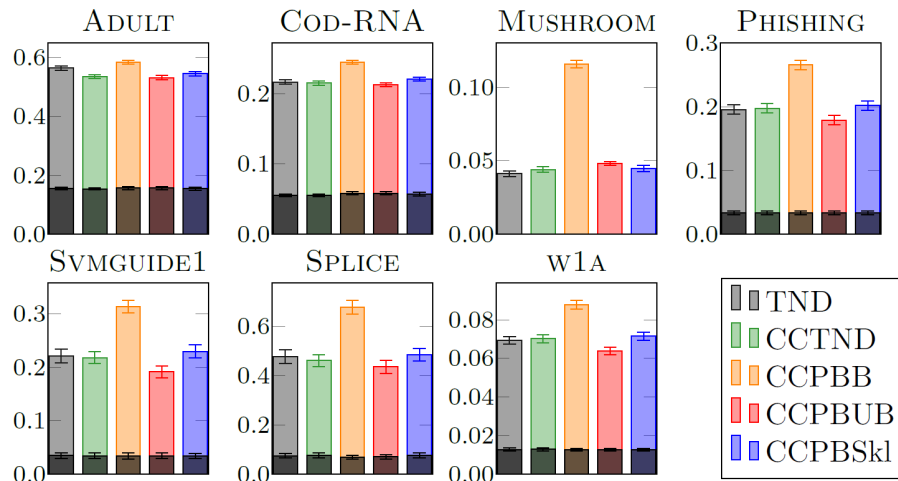


Figure 3.8: Comparison of the bounds and the test losses of the weighted majority vote on ensembles of heterogeneous classifiers with optimized posterior ρ^* generated by TND, CCTND, CCPBB, CCPBUB, and CCPBSkl. The data sets are binary labeled. The test losses of the corresponding bounds are shown in black. We report the mean and the standard deviation over 10 runs of the experiments.

and optimization methods in this experiment are the same as described in Sec. 3.6.5.3 under “Bounds and Optimization”. The source code for replicating the experiments is available at Github⁶.

Experimental Setting

In this section, we follow the construction used by Wu et al. (2021). We construct the ensemble from decision trees, which is available in *scikit-learn*. We take 100 fully grown trees to build the random forest. The ensemble is again construct by bagging as described in 3.6.5.3, where each tree h is trained on a subset of a random split S_h and estimated on \bar{S}_h . To train each tree, we use the Gini criterion for splitting and consider \sqrt{d} features in each split, where d is the number of the attributes in data.

Results

The results of random forest weighted majority vote on binary data sets are shown in Figure 3.10 while the results on multiclass data sets are shown in Figure 3.11. Similar to the discussions in Sec. 3.6.5.3 for the ensemble of heterogeneous classifiers,

⁶<https://github.com/StephanLorenzen/MajorityVoteBounds>

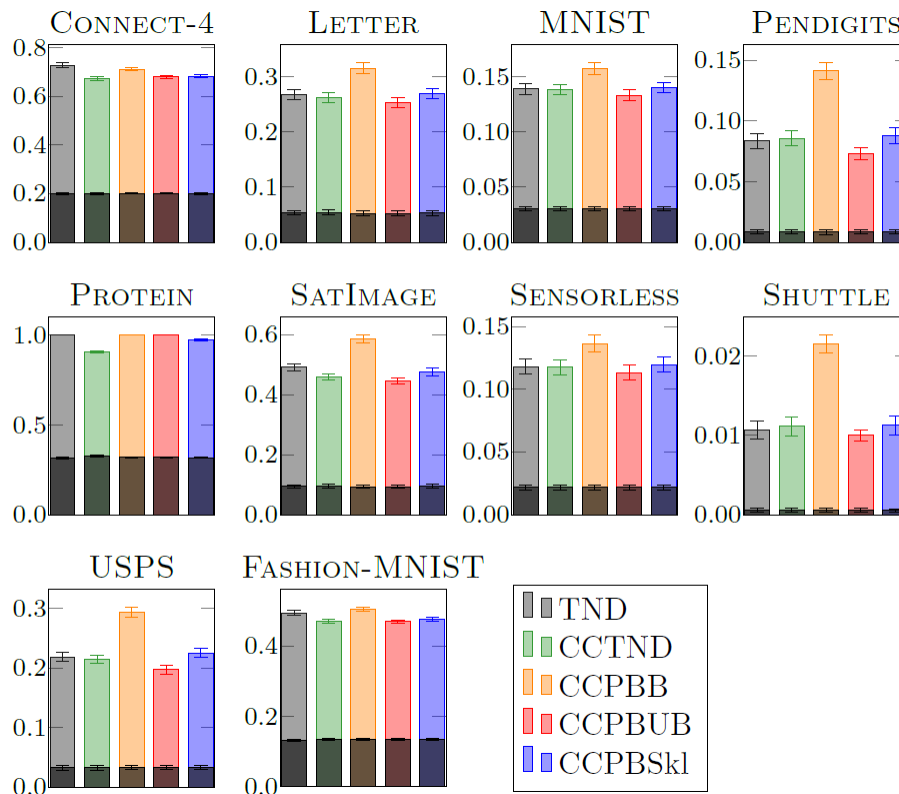


Figure 3.9: Comparison of the bounds and the test losses of the weighted majority vote on ensembles of heterogeneous classifiers with optimized posterior ρ^* generated by TND, CCTND, CCPBB, CCPBUB, and CCPBSkl. The data sets are multiclass labeled. The test losses of the corresponding bounds are shown in black. We report the mean and the standard deviation over 10 runs of the experiments.

TND serves as a baseline. In both figures, CCTND performs similar to the baseline. The CCPBB, using the α -tandem loss, lags behind due to nested application of concentration bounds. The two new bounds based the α -tandem loss, CCPBUB and CCPBSkl, clearly improve the shortage of the CCPBB bound. The CCPBSkl bound is comparable to the baseline, and the CCPBUB bound often provide tighter bounds than the baseline.

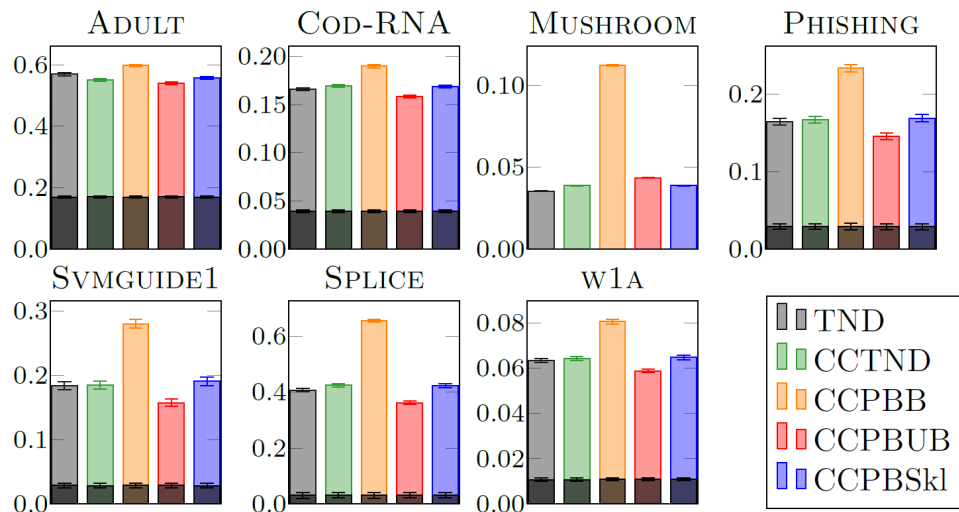


Figure 3.10: Comparison of the bounds and the test losses of the weighted majority vote on random forest with optimized posterior ρ^* generated by TND, CCTND, CCPBB, CCPBUB, and CCPBSkl. The data sets are binary labeled. The test losses of the corresponding bounds are shown in black. We report the mean and the standard deviation over 10 runs of the experiments.

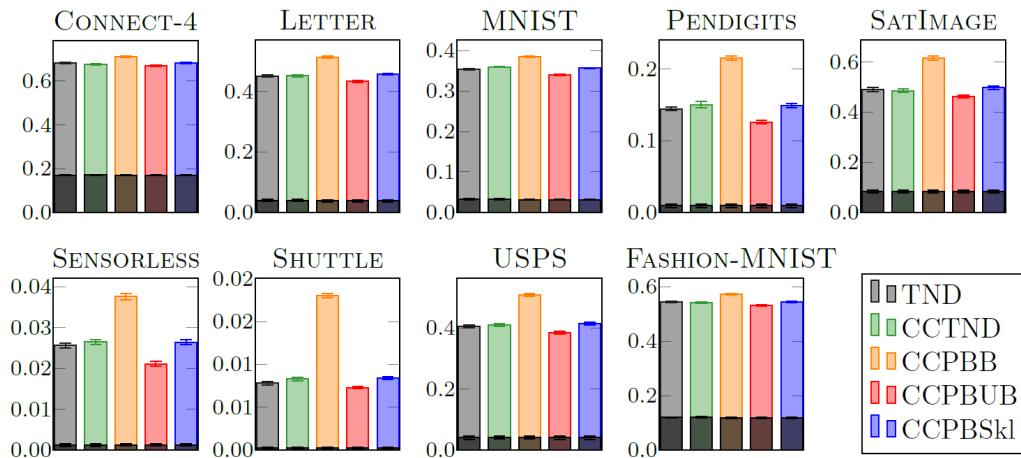


Figure 3.11: Comparison of the bounds and the test losses of the weighted majority vote on random forest with optimized posterior ρ^* generated by TND, CCTND, CCPBB, CCPBUB, and CCPBSkl. The data sets are multiclass labeled. The test losses of the corresponding bounds are shown in black. We report the mean and the standard deviation over 10 runs of the experiments.

3.6.6 PAC-Bayes- λ Inequality

Theorem 3.10 (PAC-Bayes- λ Inequality, Thiemann et al., 2017; Masegosa et al., 2020). *For any loss $\ell \in [0, 1]$, any probability distribution π on \mathcal{H} that is independent of S and any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over a random draw of a sample S , for all distributions ρ on \mathcal{H} and all $\lambda \in (0, 2)$ and $\gamma > 0$ simultaneously:*

$$\mathbb{E}_\rho [L(h)] \leq \frac{\mathbb{E}_\rho[\hat{L}(h, S)]}{1 - \frac{\lambda}{2}} + \frac{\text{KL}(\rho||\pi) + \ln(2\sqrt{n}/\delta)}{\lambda(1 - \frac{\lambda}{2})n}, \quad (3.22)$$

$$\mathbb{E}_\rho [L(h)] \geq \left(1 - \frac{\gamma}{2}\right) \mathbb{E}_\rho[\hat{L}(h, S)] - \frac{\text{KL}(\rho||\pi) + \ln(2\sqrt{n}/\delta)}{\gamma n}. \quad (3.23)$$

The upper bound is due to Thiemann et al. (2017) and the lower bound is due to Masegosa et al. (2020), and both hold simultaneously. The PAC-Bayes- λ bound is an optimization friendly relaxation of the PAC-Bayes-kl bound. Tolstikhin and Seldin (2013) has shown that for a given ρ , equation 3.22 is convex in λ and has the minimizer

$$\lambda_\rho^* = \frac{2}{\sqrt{\frac{2n\mathbb{E}_\rho[\hat{L}(h, S)]}{\text{KL}(\rho||\pi) + \ln \frac{2\sqrt{n}}{\delta}} + 1} + 1}.$$

On the other hand, Masegosa et al. (2020) has shown that for a given ρ , the optimal γ in equation 3.23 can be achieved by

$$\gamma_\rho^* = \sqrt{\frac{\text{KL}(\rho||\pi) + \ln(2\sqrt{n}/\delta)}{n\mathbb{E}_\rho[\hat{L}(h, S)]}}.$$

Furthermore, given λ or γ , the optimal ρ can be achieved by projected gradient descent.

Chapter 4

Summary and Discussion

In Chapter 2, we proposed an optimization-friendly form of the Chebyshev-Cantelli inequality and applied it to derive two new forms of second-order oracle bounds for the weighted majority vote. The new bounds bridge between the C-bounds (Germain et al., 2015) and the tandem bound (Masegosa et al., 2020) and take the best of both: the tightness of the C-bounds and the optimization and estimation convenience of the tandem bound. It would be interesting to explore other potential uses for this new version of the Chebyshev-Cantelli inequality.

In Chapter 3, we studied the concentration of measure inequalities for ternary random variables. We proposed the split-kl inequality that addresses a long-standing open question on how to exploit the structure of ternary random variables. Notably, for ternary random variables, we found that the middle value is the best choice of the parameter μ for the split-kl inequality. This choice results in the split-kl inequality being tight for ternary random variables, as the kl inequality is tight for binary random variables.

In principle, the split-kl inequality can also be applied to any bounded random variables. However, the question of how to select an appropriate value of μ for such cases remains unresolved. In our experiments with bounded random variables, we found that choosing the central value of the range may yield satisfactory results, although we do not yet have a fully satisfactory theoretical explanation for this choice.

In Chapter 3, we also included comparisons between the Unexpected Bernstein inequalities and the Empirical Bernstein inequality. The Empirical Bernstein inequality is derived by combining a Bernstein inequality for the loss with an inequality for the

variance. Due to the two-step approach, the inequality often lacks tightness. In contrast, the Unexpected Bernstein is a single-step approach that provides an empirical bound for the loss using the empirical second moment. However, its non-centered nature can lead to loose results when the second moment is significantly larger than the variance, which happens when the expectation of the random variable deviates significantly from 0. Therefore, an interesting open problem is to develop a centered version of the Unexpected Bernstein inequality that would be based on the empirical variance and still have a one-step estimation. Additionally, the Bernstein-type inequalities we presented all require a grid of the parameter γ and a union bound over the grid, such that the high-probability argument follows. It is, therefore, an interesting open question whether such a grid can be avoided.

This thesis presents the PAC-Bayes-Bennett inequality and the PAC-Bayes-Split-kl inequality for randomized and ensemble classifiers. In particular, the PAC-Bayes-Bennett inequality improves on the PAC-Bayes-Bernstein inequality proposed by Seldin et al. (2012). Meanwhile, the PAC-Bayes-Split-kl inequality retains the tightness achieved by the split-kl inequality when applied to ternary random variables. These new inequalities, along with a new version of the Chebyshev-Cantelli inequality discussed above, provide better generalization guarantees for the weighted majority vote.

List of Publications

The work presented in this thesis has lead to the following publications.

1. Yi-Shan Wu, Andres Masegosa, Stephan Sloth Lorenzen, Christian Igel, and Yevgeny Seldin. Chebyshev-cantelli pac-bayes-bennett inequality for the weighted majority vote. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
2. Yi-Shan Wu and Yevgeny Seldin. Split-kl and pac-bayes-split-kl inequalities for ternary random variables. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.

Bibliography

- Amiran Ambroladze, Emilio Parrado-Hernández, and John Shawe-Taylor. Tighter PAC-Bayes bounds. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2007.
- Jean Yves Audibert, Rémi Munos, and Csaba Szepesvári. Exploration-exploitation trade-off using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 2009.
- Baptiste Bauvin, Cécile Capponi, Jean-François Roy, and François Laviolette. Fast greedy c -bound minimization with guarantees. 2020.
- Daniel Berend and Aryeh Kontorovich. A finite sample analysis of the naive Bayes classifier. *Journal of Machine Learning Research*, 2016.
- Sergei N. Bernstein. *Probability Theory*. Moscow-Leningrad, 4th edition, 1946. In Russian.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2013.
- Leo Breiman. Bagging predictors. *Machine learning*, 24(2), 1996.
- Leo Breiman. Random forests. *Machine learning*, 45(1), 2001.
- Nicolò Cesa-Bianchi, Yishay Mansour, and Gilles Stoltz. Improved second-order bounds for prediction with expert advice. *Machine Learning*, 66, 2007.
- Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2, 2011.

- Tianqi Chen and Carlos Guestrin. XGBoost: A scalable tree boosting system. In *Proceedings of the International Conference on Knowledge Discovery and Data Mining (ACM SIGKDD)*, pages 785–794, 2016.
- Robert M Corless, Gaston H Gonnet, David EG Hare, David J Jeffrey, and Donald E Knuth. On the lambertw function. *Advances in Computational mathematics*, 5(1), 1996.
- Corinna Cortes, Giulia DeSalvo, Claudio Gentile, Mehryar Mohri, and Scott Yang. Online learning with abstention. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2018.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley Series in Telecommunications and Signal Processing, 2nd edition, 2006.
- Dheeru Dua and Casey Graff. UCI machine learning repository, 2019. URL <http://archive.ics.uci.edu/ml>.
- Gintare Karolina Dziugaite and Daniel M. Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. In *UAI*, 2017.
- Xiequan Fan, Ion Grama, and Quansheng Liu. Exponential inequalities for martingales with applications. *Electronic Journal of Probability*, 20:1–22, 2015.
- Ciprian Florescu and Christian Igel. Resilient backpropagation (rprop) for batch-learning in tensorflow. In *International Conference on Learning Representations (ICLR) Workshop*, 2018.
- Andrew Foong, Wessel Bruinsma, David Burt, and Richard Turner. How tight can pac-bayes be in the small data regime? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- Andrew Y. K. Foong, Wessel P. Bruinsma, and David R. Burt. A note on the chernoff bound for random variables in the unit interval. *arXiv preprint arXiv.2205.07880*, 2022.
- Yoav Freund and Robert E. Schapire. Experiments with a new boosting algorithm. In *Proceedings of the International Conference on Machine Learning (ICML)*, 1996.
- Jerome H. Friedman. Stochastic gradient boosting. *Computational Statistics and Data Analysis*, 38, 1999.

- Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of Statistics*, pages 1189–1232, 2001.
- Pascal Germain, Alexandre Lacasse, François Laviolette, and Mario Marchand. PAC-Bayesian learning of linear classifiers. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2009.
- Pascal Germain, Alexandre Lacasse, François Laviolette, Mario Marchand, and Jean-François Roy. Risk bounds for the majority vote: From a PAC-Bayesian analysis to a learning algorithm. *Journal of Machine Learning Research*, 16, 2015.
- W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- Christian Igel and Michael Hüsken. Empirical evaluation of the improved rprop learning algorithms. *Neurocomputing*, 50, 2003.
- Alexandre Lacasse, François Laviolette, Mario Marchand, Pascal Germain, and Nicolas Usunier. PAC-Bayes bounds for the risk of the majority vote and the variance of the Gibbs classifier. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2007.
- John Langford. Tutorial on practical prediction theory for classification. *Journal of Machine Learning Research*, 6, 2005.
- John Langford and John Shawe-Taylor. PAC-Bayes & margins. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2002.
- François Laviolette, Mario Marchand, and Jean-François Roy. From PAC-Bayes bounds to quadratic programs for majority votes. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2011.
- François Laviolette, Emilie Morvant, Liva Ralainvola, and Jean-François Roy. Risk upper bounds for general ensemble methods with an application to multiclass classification. *Neurocomputing*, 219, 2017.
- Gaël Letarte, Pascal Germain, Benjamin Guedj, and François Laviolette. Dichotomize and generalize: PAC-Bayesian binary activated deep neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Stephan Sloth Lorenzen, Christian Igel, and Yevgeny Seldin. On PAC-Bayesian bounds for random forests. *Machine Learning*, 2019.

- Andrés R. Masegosa, Stephan S. Lorenzen, Christian Igel, and Yevgeny Seldin. Second order PAC-Bayesian bounds for the weighted majority vote. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- Llew Mason, Jonathan Baxter, Peter Bartlett, and Marcus Frean. Boosting algorithms as gradient descent. In *Advances in Neural Information Processing Systems (NeurIPS)*, 1999.
- Andreas Maurer. A note on the PAC-Bayesian theorem. arXiv preprint cs/0411099, 2004.
- Andreas Maurer and Massimiliano Pontil. Empirical Bernstein bounds and sample variance penalization. In *Proceedings of the Conference on Learning Theory (COLT)*, 2009.
- David McAllester. Some PAC-Bayesian theorems. In *Proceedings of the Conference on Learning Theory (COLT)*, 1998.
- David McAllester. PAC-Bayesian stochastic model selection. *Machine Learning*, 51, 2003.
- Zakaria Mhammedi, Peter Grünwald, and Benjamin Guedj. PAC-Bayes un-expected Bernstein inequality. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- Volodymyr Mnih, Csaba Szepesvári, and Jean-Yves Audibert. Empirical Bernstein stopping. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2008.
- David Reeb, Andreas Doerr, Sebastian Gerwin, and Barbara Rakitsch. Learning gaussian processes by minimizing PAC-Bayesian generalization bounds. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Matthias Seeger. PAC-Bayesian generalization error bounds for Gaussian process classification. *Journal of Machine Learning Research*, 3, 2002.
- Yevgeny Seldin, François Laviolette, Nicolò Cesa-Bianchi, John Shawe-Taylor, and Peter Auer. PAC-Bayesian inequalities for martingales. *IEEE Transactions on Information Theory*, 58, 2012.
- Niklas Thiemann, Christian Igel, and Yevgeny Seldin. PAC-Bayesian aggregation without cross-validation. <http://arxiv.org/abs/1608.05610>, 2016.

- Niklas Thiemann, Christian Igel, Olivier Wintenberger, and Yevgeny Seldin. A strongly quasiconvex PAC-Bayesian bound. In *Proceedings of the International Conference on Algorithmic Learning Theory (ALT)*, 2017.
- Sunil Thulasidasan, Tanmoy Bhattacharya, Jeff A. Bilmes, Gopinath Chennupati, and Jamal Mohd-Yusof. Combating label noise in deep learning using abstention. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.
- Ilya Tolstikhin and Yevgeny Seldin. PAC-Bayes-Empirical-Bernstein inequality. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2013.
- Paul Viallard, Pascal Germain, Amaury Habrard, and Emilie Morvant. Self-bounding majority vote learning algorithms by the direct minimization of a tight PAC-Bayesian C-bound. In *Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)*, 2021.
- Olivier Wintenberger. Optimal learning with Bernstein online aggregation. *Machine Learning*, 106, 2017.
- Yi-Shan Wu and Yevgeny Seldin. Split-kl and pac-bayes-split-kl inequalities for ternary random variables. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- Yi-Shan Wu, Andres Masegosa, Stephan Sloth Lorenzen, Christian Igel, and Yevgeny Seldin. Chebyshev-cantelli pac-bayes-bennett inequality for the weighted majority vote. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- Wenda Zhou, Victor Veitch, Morgane Austern, Ryan P. Adams, and Peter Orbanz. Non-vacuous generalization bounds at the imagenet scale: a PAC-bayesian compression approach. In *ICLR*, 2019.