

Controlling Confusion via Generalisation Bounds

Anonymous authors

Paper under double-blind review

Abstract

We establish new generalisation bounds for multiclass classification by abstracting to a more general setting of discretised error types. Extending the PAC-Bayes theory, we are hence able to provide fine-grained bounds on performance for multiclass classification, as well as applications to other learning problems including discretisation of regression losses. Tractable training objectives are derived from the bounds. The bounds are uniform over all weightings of the discretised error types and thus can be used to bound weightings not foreseen at training, including the full confusion matrix in the multiclass classification case.

1 Introduction

Generalisation bounds are a core component of the theoretical understanding of machine learning algorithms. For over two decades now, the PAC-Bayesian theory has been at the core of studies on generalisation abilities of machine learning algorithms. PAC-Bayes originates in the seminal work of McAllester (1998; 1999) and was further developed by Catoni (2003; 2004; 2007), among other authors—we refer to the recent surveys Guedj (2019) and Alquier (2021) for an introduction to the field. The outstanding empirical successes of deep neural networks in the past decade call for better theoretical understanding of deep learning, and PAC-Bayes emerged as one of the few frameworks allowing the derivation of meaningful (and non-vacuous) generalisation bounds for neural networks: the pioneering work of Dziugaite & Roy (2017) has been followed by a number of contributions, including Neyshabur et al. (2018), Zhou et al. (2019), Letarte et al. (2019), Pérez-Ortiz et al. (2021); Pérez-Ortiz et al. (2021) and Biggs & Guedj (2021; 2022a;b), to name but a few.

Much of the PAC-Bayes literature focuses on the case of binary classification, or of multiclass classification where one only distinguishes whether each classification is correct or incorrect. This is in stark contrast to the complexity of contemporary real-world learning problems. This work aims to bridge this gap via generalisation bounds that provide information rich measures of performance at test time by controlling the probabilities of errors of any finite number of types, bounding combinations of these probabilities uniformly over all weightings.

Previous results. We believe our framework of discretised error types to be novel. In the particular case of multiclass classification, little is known from a theoretical perspective and, to the best of our knowledge, only a handful of relevant strategies or generalisation bounds can be compared to the present paper. The closest is the work of Morvant et al. (2012) on a PAC-Bayes generalisation bound on the operator norm of the confusion matrix, to train a Gibbs classifier. We focus on a different performance metric, in the broader setting of discretised error types. Koço & Capponi (2013) suggest to minimise the confusion matrix norm with a focus on the imbalance between classes; their treatment is not done through PAC-Bayes. Laviolette et al. (2017) extend the celebrated \mathcal{C} -bound in PAC-Bayes to weighted majority votes of classifiers, to perform multiclass classification. Benabbou & Lang (2017) present a streamlined version of some of the results from Morvant et al. (2012) in the case where some examples are voluntarily not classified (*e.g.*, in the case of too large uncertainty). More recently, Feofanov et al. (2019) derive bounds for a majority vote classifier where the confusion matrix serves as an error indicator: they conduct a study of the Bayes classifier.

From binary to multiclass classification. A number of PAC-Bayesian bounds have been unified by a single general bound, found in Bégin et al. (2016). Stated as Theorem 1 below, it applies to binary classification. We use it as a basis to prove our Theorem 3, a more general bound that can be applied to,

amongst other things, multiclass classification and discretised regression. While the proof of Theorem 3 follows similar lines to that given in Bégin et al. (2016), our generalisation to ‘soft’ hypotheses incurring any finite number of error types requires a non-trivial extension of a result found in Maurer (2004). This extension (Lemma 5), along with its corollary (Corollary 6) may be of independent interest. The generalisation bound in Maurer (2004), stated below as Corollary 2, is shown in Bégin et al. (2016) to be a corollary of their bound. In a similar manner, we derive Corollary 7 from Theorem 3. Obtaining this corollary is significantly more involved than the analogous derivation in Bégin et al. (2016) or the original proof in Maurer (2004), requiring a number of technical results found in Appendix B.

Briefly, the results in Bégin et al. (2016) and Maurer (2004) consider an arbitrary input set \mathcal{X} , output set $\mathcal{Y} = \{-1, 1\}$, hypothesis space $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ and i.i.d. sample $S \in (\mathcal{X} \times \mathcal{Y})^m$. They then establish high probability bounds on the discrepancy between the risk (probability of error on a new datapoint) of any stochastic classifier Q (namely, a distribution on \mathcal{H}) and its empirical counterpart (the fraction of the sample Q misclassifies). The bounds hold uniformly over all Q and contain a complexity term involving the Kullback-Leibler (KL) divergence between Q and a reference distribution P on \mathcal{H} (often referred to as a prior by analogy with Bayesian inference—see the discussion in Guedj, 2019).

There are two ways in which the results in Bégin et al. (2016) and Maurer (2004) can be described as binary. First, as \mathcal{Y} contains two elements, this is obviously an instance of binary classification. But a more interesting and subtle way to look at this is that only two cases are distinguished—correct classification and incorrect classification. Specifically, since the two different directions in which misclassification can be made are counted together, the bound gives no information on which direction is more likely.

More generally, the aforementioned bounds can be applied in the context of multiclass classification provided one maintains the second binary characteristic by only distinguishing correct and incorrect classifications rather than considering the entire confusion matrix. However, note that these bounds will not give information on the relative likelihood of the different errors. In contrast, our new results can consider the entire confusion matrix, bounding how far the true (read “expected over the data-generating distribution”) confusion matrix differs from the empirical one, according to some metric. In fact, our results extend to the case of arbitrary label set \mathcal{Y} , provided the number of different errors one distinguishes is finite.

Formally, we let $\bigcup_{j=1}^M E_j$ be a user-specified disjoint partition of \mathcal{Y}^2 into a finite number of M error types, where we say that a hypothesis $h \in \mathcal{H}$ makes an error of type j on datapoint (x, y) if $(h(x), y) \in E_j$ (by convention, every pair $(\hat{y}, y) \in \mathcal{Y}^2$ is interpreted as a predicted value \hat{y} followed by a true value y , in that order). It should be stressed that some E_j need not correspond to mislabellings—indeed, some of the E_j may distinguish different correct labellings. We then count up the number of errors of each type that a hypothesis makes on a sample, and bound how far this empirical distribution of errors is from the expected distribution under the data-generating distribution (Theorem 3). Thus, in our generalisation, the (scalar) risk and empirical risk ($R_D(Q)$ and $R_S(Q)$, defined in the next section) are replaced by M -dimensional vectors ($\mathbf{R}_D(Q)$ and $\mathbf{R}_S(Q)$), and our discrepancy measure d is a divergence between discrete distributions on M elements. Our generalisation therefore allows us to bound how far the true distribution of errors can be from the observed distribution of errors. If we then associate a loss value $\ell_j \in [0, \infty)$ to each E_j we can derive a bound on the *total risk*, defined as the sum of the true error probabilities weighted by the loss values. In fact, the total risk is bounded with high probability uniformly over all such weightings. The loss values need not be distinct; we may wish to understand the distribution of error types even across error types that incur the same loss.

For example, in the case of binary classification with $\mathcal{Y} = \{-1, 1\}$, we can take the usual partition into $E_1 = \{(-1, -1), (1, 1)\}$ and $E_2 = \{(-1, 1), (1, -1)\}$ and loss values $\ell_1 = 0, \ell_2 = 1$, or the fine-grained partition $\mathcal{Y}^2 = \{(0, 0)\} \cup \{(1, 1)\} \cup \{(0, 1)\} \cup \{(1, 0)\}$ and the loss values $\ell_1 = \ell_2 = 0, \ell_3 = 1, \ell_4 = 2$. More generally, for multiclass classification with N classes and $\mathcal{Y} = [N]$, one may take the usual coarse partition into $E_1 = \{(\hat{y}, y) \in \mathcal{Y}^2 : \hat{y} = y\}$ and $E_2 = \{(\hat{y}, y) \in \mathcal{Y}^2 : \hat{y} \neq y\}$ (with $\ell_1 = 0$ and $\ell_2 = 1$), or the fully refined partition into $E_{i,j} = \{(i, j)\}$ for $i, j \in [N]$ (with correspondingly greater choice of the associated loss values), or something in-between. Note that we still refer to E_j as an “error type” even if it contains elements that correspond to correct classification, namely if there exists $y \in \mathcal{Y}$ such that $(y, y) \in E_j$. As we will see later, a more fine-grained partition will allow more error types to be distinguished and bounded, at the expense of a

looser bound. As a final example, for regression with $\mathcal{Y} = \mathbb{R}$, we may fix M strictly increasing thresholds $0 = \lambda_1 < \lambda_2 < \dots < \lambda_M$ and partition \mathcal{Y}^2 into $E_j = \{(y_1, y_2) \in \mathcal{Y}^2 : \lambda_j \leq |y_1 - y_2| < \lambda_{j+1}\}$ for $j \in [M-1]$, and $E_M = \{(y_1, y_2) \in \mathcal{Y}^2 : |y_1 - y_2| \geq \lambda_M\}$.

Outline. We set our notation in Section 2. In Section 3 we state and prove generalisation bounds in the setting of discretised error types: this significantly expands the previously known results from Bégin et al. (2016) by allowing for generic output sets \mathcal{Y} . Our main results are Theorem 3 and Corollary 7. To make our findings profitable to the broader machine learning community we then discuss how these new bounds can be turned into tractable training objectives in Section 4 (with a general recipe described in greater detail in Appendix A). The paper closes with perspectives for follow-up work in Section 5 and we defer to Appendix B the proofs of technical results.

2 Notation

For any set A , let $\mathcal{M}(A)$ be the set of probability measures on A . For any $M \in \mathbb{Z}_{>0}$, define $[M] := \{1, 2, \dots, M\}$, the M -dimensional simplex $\Delta_M := \{\mathbf{u} \in [0, 1]^M : u_1 + \dots + u_M = 1\}$ and its interior $\Delta_M^{\circ} := \Delta_M \cap (0, 1)^M$. For $m, M \in \mathbb{Z}_{>0}$, define the integer counterparts $S_{m, M} := \{(k_1, \dots, k_M) \in \mathbb{Z}_{\geq 0}^M : k_1 + \dots + k_M = m\}$ and $S_{m, M}^{\circ} := S_{m, M} \cap \mathbb{Z}_{>0}^M$. The set $S_{m, M}$ is the domain of the multinomial distribution with parameters m, M and some $\mathbf{r} \in \Delta_M$, which is denoted $\text{Mult}(m, M, \mathbf{r})$ and has probability mass function for $\mathbf{k} \in S_{m, M}$ given by

$$\text{Mult}(\mathbf{k}; m, M, \mathbf{r}) := \binom{m}{k_1 \ k_2 \ \dots \ k_M} \prod_{j=1}^M r_j^{k_j}, \quad \text{where} \quad \binom{m}{k_1 \ k_2 \ \dots \ k_M} := \frac{m!}{\prod_{j=1}^M k_j!}.$$

For $\mathbf{q}, \mathbf{p} \in \Delta_M$, let $\text{kl}(\mathbf{q} \parallel \mathbf{p})$ denote the KL-divergence of $\text{Mult}(1, M, \mathbf{q})$ from $\text{Mult}(1, M, \mathbf{p})$, namely $\text{kl}(\mathbf{q} \parallel \mathbf{p}) := \sum_{j=1}^M q_j \ln \frac{q_j}{p_j}$, with the convention that $0 \ln \frac{0}{x} = 0$ for $x \geq 0$ and $x \ln \frac{x}{0} = \infty$ for $x > 0$. For $M = 2$ we abuse notation and abbreviate $\text{kl}((q, 1-q) \parallel (p, 1-p))$ to $\text{kl}(q \parallel p)$, which is then the conventional definition of $\text{kl}(\cdot \parallel \cdot) : [0, 1]^2 \rightarrow [0, \infty]$ found in the PAC-Bayes literature (as in Seeger, 2002, for example).

Let \mathcal{X} and \mathcal{Y} be arbitrary input (*e.g.*, feature) and output (*e.g.*, label) sets respectively. Let $\bigcup_{j=1}^M E_j$ be a partition of \mathcal{Y}^2 into a finite sequence of M error types, and to each E_j associate a loss value $\ell_j \in [0, \infty)$. The only restriction we place on the loss values ℓ_j is that they are not all equal. This is not a strong assumption, since if they were all equal then all hypotheses would incur equal loss and there would be no learning problem: we are effectively ruling out trivial cases.

Let $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ denote a hypothesis class, $D \in \mathcal{M}(\mathcal{X} \times \mathcal{Y})$ a data-generating distribution and $S \sim D^m$ an i.i.d. sample of size m drawn from D . For $h \in \mathcal{H}$ and $j \in [M]$ we define the *empirical j -risk* and *true j -risk* of h to be $R_S^j(h) := \frac{1}{m} \sum_{(x, y) \in S} \mathbb{1}[(h(x), y) \in E_j]$ and $R_D^j(h) := \mathbb{E}_{(x, y) \sim D} [\mathbb{1}[(h(x), y) \in E_j]]$, respectively, namely, the proportion of the sample S on which h makes an error of type E_j and the probability that h makes an error of type E_j on a new $(x, y) \sim D$.

More generally, suppose $\mathcal{H} \subseteq \mathcal{M}(\mathcal{Y})^{\mathcal{X}}$ is a class of *soft* hypotheses of the form $H : \mathcal{X} \rightarrow \mathcal{M}(\mathcal{Y})$, where, for any $A \subseteq \mathcal{Y}$, $H(x)[A]$ is interpreted as the probability according to H that the label of x is in A . It is worth stressing that a soft hypothesis is still deterministic since a prediction is not drawn from the distribution it returns. We then define the *empirical j -risk* of H to be $R_S^j(H) := \frac{1}{m} \sum_{(x, y) \in S} H(x)[\{\hat{y} \in \mathcal{Y} : (\hat{y}, y) \in E_j\}]$, namely the mean—over the elements (x, y) of S —probability mass H assigns to predictions $\hat{y} \in \mathcal{Y}$ incurring an error of type E_j when labelling each x . Further, we define the *true j -risk* of H to be $R_D^j(H) := \mathbb{E}_{(x, y) \sim D} [H(x)[\{\hat{y} \in \mathcal{Y} : (\hat{y}, y) \in E_j\}]]$, namely the mean—over $(x, y) \sim D$ —probability mass H assigns to predictions $\hat{y} \in \mathcal{Y}$ incurring an error of type E_j when labelling each x . We will see in Section 4 that the more general hypothesis class $\mathcal{H} \subseteq \mathcal{M}(\mathcal{Y})^{\mathcal{X}}$ is necessary for constructing a differentiable training objective.

To each ordinary hypothesis $h \in \mathcal{Y}^{\mathcal{X}}$ there corresponds a soft hypothesis $H \in \mathcal{M}(\mathcal{Y})^{\mathcal{X}}$ that, for each $x \in \mathcal{X}$, returns a point mass on $h(x)$. In this case, it is straightforward to show that $R_S^j(h) = R_S^j(H)$ and $R_D^j(h) = R_D^j(H)$ for all $j \in [M]$, where we have used the corresponding definitions above for ordinary and

soft hypotheses. Since, in addition, our results hold identically for both ordinary and soft hypotheses, we henceforth use the same notation h for both ordinary and soft hypotheses and their associated values $R_S^j(h)$ and $R_D^j(h)$. It will always be clear from the context whether we are dealing with ordinary or soft hypotheses and thus which of the above definitions of the empirical and true j -risks is being used.

We define the *empirical risk* and *true risk* of a (ordinary or soft) hypothesis h to be $\mathbf{R}_S(h) := (R_S^1(h), \dots, R_S^M(h))$ and $\mathbf{R}_D(h) := (R_D^1(h), \dots, R_D^M(h))$, respectively. It is straightforward to show that $\mathbf{R}_S(h)$ and $\mathbf{R}_D(h)$ are elements of Δ_M . Since S is drawn i.i.d. from D , the expectation of the empirical risk is equal to the true risk, namely $\mathbb{E}_S[R_S^j(h)] = R_D^j(h)$ for all j and thus $\mathbb{E}_S[\mathbf{R}_S(h)] = \mathbf{R}_D(h)$. Finally, we generalise to stochastic hypotheses $Q \in \mathcal{M}(\mathcal{H})$, which predict by first drawing a deterministic hypothesis $h \sim Q$ and then predicting according to h , where a new h is drawn for each prediction. Thus, we define the *empirical j -risk* and *true j -risk* of Q to be the scalars $R_S^j(Q) := \mathbb{E}_{h \sim Q}[R_S^j(h)]$ and $R_D^j(Q) := \mathbb{E}_{h \sim Q}[R_D^j(h)]$, for $j \in [M]$, and simply the *empirical risk* and *true risk* of Q to be the elements of Δ_M defined by $\mathbf{R}_S(Q) := \mathbb{E}_{h \sim Q}[\mathbf{R}_S(h)]$ and $\mathbf{R}_D(Q) := \mathbb{E}_{h \sim Q}[\mathbf{R}_D(h)]$. As before, since S is i.i.d., we have (using Fubini this time) that $\mathbb{E}_S[\mathbf{R}_S(Q)] = \mathbf{R}_D(Q)$. Finally, given a loss vector $\ell \in [0, \infty)^M$, we define the *total risk* of Q by the scalar $R_D^T(Q) := \sum_{j=1}^M \ell_j R_D^j(Q)$. As is conventional in the PAC-Bayes literature, we refer to sample independent and dependent distributions on $\mathcal{M}(\mathcal{H})$ (*i.e.* stochastic hypotheses) as *priors* (denoted P) and *posteriors* (denoted Q) respectively, even if they are not related by Bayes' theorem.

3 Inspiration and Main Results

We first state the existing results in [Bégin et al. \(2016\)](#) and [Maurer \(2004\)](#) that we will generalise from just two error types (correct and incorrect) to any finite number of error types. These results are stated in terms of the scalars $R_S(Q) := \frac{1}{m} \sum_{(x,y) \in S} \mathbb{1}[h(x) \neq y]$ and $R_D(Q) := \mathbb{E}_{(x,y) \sim D} \mathbb{1}[h(x) \neq y]$ and, as we demonstrate, correspond to the case $M = 2$ of our generalisations.

Theorem 1. ([Bégin et al., 2016, Theorem 4](#)) *Let \mathcal{X} be an arbitrary set and $\mathcal{Y} = \{-1, 1\}$. Let $D \in \mathcal{M}(\mathcal{X} \times \mathcal{Y})$ be a data-generating distribution and $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ be a hypothesis class. For any prior $P \in \mathcal{M}(\mathcal{H})$, $\delta \in (0, 1]$, convex function $d : [0, 1]^2 \rightarrow \mathbb{R}$, sample size m and $\beta \in (0, \infty)$, with probability at least $1 - \delta$ over the random draw $S \sim D^m$, we have that simultaneously for all posteriors $Q \in \mathcal{M}(\mathcal{H})$*

$$d(R_S(Q), R_D(Q)) \leq \frac{1}{\beta} \left[\text{KL}(Q \| P) + \ln \frac{\mathcal{I}_d(m, \beta)}{\delta} \right],$$

with $\mathcal{I}_d(m, \beta) := \sup_{r \in [0, 1]} \left[\sum_{k=0}^m \text{Bin}(k; m, r) \exp \left(\beta d \left(\frac{k}{m}, r \right) \right) \right]$, where $\text{Bin}(k; m, r)$ is the binomial probability mass function $\text{Bin}(k; m, r) := \binom{m}{k} r^k (1-r)^{m-k}$.

Note the original statement in [Bégin et al. \(2016\)](#) is for a positive integer m' , but the proof trivially generalises to any $\beta \in (0, \infty)$. One of the bounds that Theorem 1 unifies—which we also generalise—is that of [Seeger \(2002\)](#), later tightened in [Maurer \(2004\)](#), which we now state. It can be recovered from Theorem 1 by setting $\beta = m$ and $d(q, p) = \text{kl}(q \| p) := q \ln \frac{q}{p} + (1-q) \ln \frac{1-q}{1-p}$.

Corollary 2. ([Maurer, 2004, Theorem 5](#)) *Let \mathcal{X} be an arbitrary set and $\mathcal{Y} = \{-1, 1\}$. Let $D \in \mathcal{M}(\mathcal{X} \times \mathcal{Y})$ be a data-generating distribution and $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ be a hypothesis class. For any prior $P \in \mathcal{M}(\mathcal{H})$, $\delta \in (0, 1]$ and sample size m , with probability at least $1 - \delta$ over the random draw $S \sim D^m$, we have that simultaneously for all posteriors $Q \in \mathcal{M}(\mathcal{H})$*

$$\text{kl}(R_S(Q), R_D(Q)) \leq \frac{1}{m} \left[\text{KL}(Q \| P) + \ln \frac{2\sqrt{m}}{\delta} \right].$$

We wish to bound the deviation of the empirical vector $\mathbf{R}_S(Q)$ from the unknown vector $\mathbf{R}_D(Q)$. Since in general the stochastic hypothesis Q we learn will depend on the sample S , it is useful to obtain bounds on the deviation of $\mathbf{R}_S(Q)$ from $\mathbf{R}_D(Q)$ that are uniform over Q , just as in Theorem 1 and Corollary 2. In Theorem 1, the deviation $d(R_S(Q), R_D(Q))$ between the scalars $R_S(Q), R_D(Q) \in [0, 1]$ is measured by some convex function $d : [0, 1]^2 \rightarrow \mathbb{R}$. In our case, the deviation $d(\mathbf{R}_S(Q), \mathbf{R}_D(Q))$ between the vectors $\mathbf{R}_S(Q), \mathbf{R}_D(Q) \in \Delta_M$ is measured by some convex function $d : \Delta_M^2 \rightarrow \mathbb{R}$. In Section 3.2 we will derive

Corollary 7 from Theorem 3 by selecting $\beta = m$ and $d(\mathbf{q}, \mathbf{p}) := \text{kl}(\mathbf{q} \parallel \mathbf{p})$, analogous to how Corollary 2 is obtained from Theorem 1.

3.1 Statement and proof of the generalised bound

We now state and prove our generalisation of Theorem 1. The proof follows identical lines to that of Theorem 1 given in Bégin et al. (2016), but with additional non-trivial steps to account for the greater number of error types and the possibility of soft hypotheses.

Theorem 3. *Let \mathcal{X} and \mathcal{Y} be arbitrary sets and $\bigcup_{j=1}^M E_j$ be a disjoint partition of \mathcal{Y}^2 . Let $D \in \mathcal{M}(\mathcal{X} \times \mathcal{Y})$ be a data-generating distribution and $\mathcal{H} \subseteq \mathcal{M}(\mathcal{Y})^{\mathcal{X}}$ be a hypothesis class. For any prior $P \in \mathcal{M}(\mathcal{H})$, $\delta \in (0, 1]$, jointly convex function $d : \Delta_M^2 \rightarrow \mathbb{R}$, sample size m and $\beta \in (0, \infty)$, with probability at least $1 - \delta$ over the random draw $S \sim D^m$, we have that simultaneously for all posteriors $Q \in \mathcal{M}(\mathcal{H})$*

$$d(\mathbf{R}_S(Q), \mathbf{R}_D(Q)) \leq \frac{1}{\beta} \left[\text{KL}(Q \parallel P) + \ln \frac{\mathcal{I}_d(m, \beta)}{\delta} \right], \quad (1)$$

where $\mathcal{I}_d(m, \beta) := \sup_{\mathbf{r} \in \Delta_M} \left[\sum_{\mathbf{k} \in S_{m, M}} \text{Mult}(\mathbf{k}; m, M, \mathbf{r}) \exp\left(\beta d\left(\frac{\mathbf{k}}{m}, \mathbf{r}\right)\right) \right]$. Further, the bounds are unchanged if one restricts to an ordinary hypothesis class, namely if $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$.

The proof begins on the following page after a discussion and some auxiliary results. One can derive multiple bounds from this theorem, all of which then hold simultaneously with probability at least $1 - \delta$. For example, one can derive bounds on the individual error probabilities $R_D^j(Q)$ or combinations thereof. It is this flexibility that allows Theorem 3 to provide far richer information on the performance of the posterior Q on unseen data. For a more in depth discussion of how such bounds can be derived, including a recipe for transforming the bound into a differentiable training objective, see Section 4 and Appendix A.

To see that Theorem 3 is a generalisation of Theorem 1, note that we can recover it by setting $\mathcal{Y} = \{-1, 1\}$, $M = 2$, $E_1 = \{(-y, y) : y \in \mathcal{Y}\}$ and $E_2 = \{(y, y) : y \in \mathcal{Y}\}$. Then, for any convex function $d : [0, 1]^2 \rightarrow \mathbb{R}$, apply Theorem 3 with the convex function $d' : \Delta_2^2 \rightarrow \mathbb{R}$ defined by $d'((u_1, u_2), (v_1, v_2)) := d(u_1, v_1)$ so that Theorem 3 bounds $d'(\mathbf{R}_S(Q), \mathbf{R}_D(Q)) = d(R_S^1(Q), R_D^1(Q))$ which equals $d(R_S(Q), R_D(Q))$ in the notation of Theorem 1. Further,

$$\sum_{\mathbf{k} \in S_{m, 2}} \text{Mult}(\mathbf{k}; m, 2, \mathbf{r}) \exp\left(\beta d'\left(\frac{\mathbf{k}}{m}, \mathbf{r}\right)\right) = \sum_{k=0}^m \text{Bin}(k; m, r_1) \exp\left(\beta d\left(\frac{k}{m}, r_1\right)\right),$$

so that the supremum over $r_1 \in [0, 1]$ of the right hand side equals the supremum over $\mathbf{r} \in \Delta_2$ of the left hand side, which, when substituted into (1), yields the bound given in Theorem 1.

Our proof of Theorem 3 follows the lines of the proof of Theorem 1 in Bégin et al. (2016), making use of the change of measure inequality Lemma 4. However, a complication arises from the use of soft classifiers $h \in \mathcal{M}(\mathcal{Y})^{\mathcal{X}}$. A similar problem is dealt with in Maurer (2004) when proving Corollary 2 by means of a Lemma permitting the replacement of $[0, 1]$ -valued random variables by corresponding $\{0, 1\}$ -valued random variables with the same mean. We use a generalisation of this, stated as Lemma 5 (Lemma 3 in Maurer, 2004 corresponds to the case $M = 2$), the proof of which is not insightful for our purposes and thus deferred to Appendix B.1. An immediate consequence of Lemma 5 is Corollary 6, which is a generalisation of the first half of Theorem 1 in Maurer (2004). While we only use it implicitly in the remainder of the paper, we state it as it may be of broader interest.

The consequence of Lemma 5 is that the worst case (in terms of bounding $d(\mathbf{R}_S(Q), \mathbf{R}_D(Q))$) occurs when $\mathbf{R}_{\{(x, y)\}}(h)$ is a one-hot vector for all $(x, y) \in S$ and $h \in \mathcal{H}$, namely when $\mathcal{H} \subseteq \mathcal{M}(\mathcal{Y})^{\mathcal{X}}$ only contains hypotheses that, when labelling S , put all their mass on elements $\hat{y} \in \mathcal{Y}$ that incur the same error type¹. In particular, this is the case for hypotheses that put all their mass on a single element of \mathcal{Y} , equivalent to the simpler case $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ as discussed in Section 2. Thus, Lemma 5 shows that the bound given in Theorem 3 cannot be made tighter only by restricting to such hypotheses.

¹More precisely, when $\forall h \in \mathcal{H} \forall (x, y) \in S \exists j \in [M]$ such that $h(x)[\{\hat{y} \in \mathcal{Y} : (\hat{y}, y) \in E_j\}] = 1$.

Lemma 4. (Change of measure, [Csiszár, 1975](#), [Donsker & Varadhan, 1975](#)) For any set \mathcal{H} , any $P, Q \in \mathcal{M}(\mathcal{H})$ and any measurable function $\phi : \mathcal{H} \rightarrow \mathbb{R}$, $\mathbb{E}_{h \sim Q} \phi(h) \leq \text{KL}(Q \| P) + \ln \mathbb{E}_{h \sim P} \exp(\phi(h))$.

Lemma 5. (Generalisation of Lemma 3 in [Maurer, 2004](#)) Let $\mathbf{X}_1, \dots, \mathbf{X}_m$ be i.i.d Δ_M -valued random vectors with mean $\boldsymbol{\mu}$ and suppose that $f : \Delta_M^m \rightarrow \mathbb{R}$ is convex. If $\mathbf{X}'_1, \dots, \mathbf{X}'_m$ are i.i.d. $\text{Mult}(1, M, \boldsymbol{\mu})$ random vectors, then $\mathbb{E}[f(\mathbf{X}_1, \dots, \mathbf{X}_m)] \leq \mathbb{E}[f(\mathbf{X}'_1, \dots, \mathbf{X}'_m)]$.

Corollary 6. (Generalisation of Theorem 1 in [Maurer, 2004](#)) Let $\mathbf{X}_1, \dots, \mathbf{X}_m$ be i.i.d Δ_M -valued random vectors with mean $\boldsymbol{\mu}$, and $\mathbf{X}'_1, \dots, \mathbf{X}'_m$ be i.i.d. $\text{Mult}(1, M, \boldsymbol{\mu})$. Define $\bar{\mathbf{X}} := \frac{1}{m} \sum_{i=1}^m \mathbf{X}_i$ and $\bar{\mathbf{X}}' := \frac{1}{m} \sum_{i=1}^m \mathbf{X}'_i$. Then $\mathbb{E}[\exp(\text{mkl}(\bar{\mathbf{X}} \| \boldsymbol{\mu}))] \leq \mathbb{E}[\exp(\text{mkl}(\bar{\mathbf{X}}' \| \boldsymbol{\mu}))]$.

Proof. (of Corollary 6) This is immediate from Lemma 5 since the average is linear, the kl-divergence is convex and the exponential is non-decreasing and convex. \square

Proof. (of Theorem 3) The case $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ follows directly from the more general case by taking $\mathcal{H}' := \{h' \in \mathcal{M}(\mathcal{Y})^{\mathcal{X}} : \exists h \in \mathcal{H} \text{ such that } \forall x \in \mathcal{X} \ h'(x) = \delta_{h(x)}\}$, where $\delta_{h(x)} \in \mathcal{M}(\mathcal{Y})$ denotes a point mass on $h(x)$. For the general case $\mathcal{H} \subseteq \mathcal{M}(\mathcal{Y})^{\mathcal{X}}$, using Jensen's inequality with the convex function $d(\cdot, \cdot)$ and Lemma 4 with $\phi(h) = \beta d(\mathbf{R}_S(h), \mathbf{R}_D(h))$, we see that for all $Q \in \mathcal{M}(\mathcal{H})$

$$\begin{aligned} \beta d(\mathbf{R}_S(Q), \mathbf{R}_D(Q)) &= \beta d\left(\mathbb{E}_{h \sim Q} \mathbf{R}_S(h), \mathbb{E}_{h \sim Q} \mathbf{R}_D(h)\right) \\ &\leq \mathbb{E}_{h \sim Q} \beta d(\mathbf{R}_S(h), \mathbf{R}_D(h)) \\ &\leq \text{KL}(Q \| P) + \ln \left(\mathbb{E}_{h \sim P} \exp\left(\beta d(\mathbf{R}_S(h), \mathbf{R}_D(h))\right) \right) \\ &= \text{KL}(Q \| P) + \ln(Z_P(S)), \end{aligned}$$

where $Z_P(S) := \mathbb{E}_{h \sim P} \exp(\beta d(\mathbf{R}_S(h), \mathbf{R}_D(h)))$. Note that $Z_P(S)$ is a non-negative random variable, so that by Markov's inequality $\mathbb{P}_{S \sim D^m} \left(Z_P(S) \leq \frac{\mathbb{E}_{S' \sim D^m} Z_P(S')}{\delta} \right) \geq 1 - \delta$. Thus, since $\ln(\cdot)$ is strictly increasing, with probability at least $1 - \delta$ over $S \sim D^m$, we have that simultaneously for all $Q \in \mathcal{M}(\mathcal{H})$

$$\beta d(\mathbf{R}_S(Q), \mathbf{R}_D(Q)) \leq \text{KL}(Q \| P) + \ln \frac{\mathbb{E}_{S' \sim D^m} Z_P(S')}{\delta}. \quad (2)$$

To bound $\mathbb{E}_{S' \sim D^m} Z_P(S')$, let $\mathbf{X}_i := \mathbf{R}_{\{(x_i, y_i)'\}}(h) \in \Delta_M$ for $i \in [m]$, where $(x_i, y_i)'$ is the i 'th element of the dummy sample S' . Noting that each \mathbf{X}_i has mean $\mathbf{R}_D(h)$, define the random vectors $\mathbf{X}'_i \sim \text{Mult}(1, M, \mathbf{R}_D(h))$ and $\mathbf{Y} := \sum_{i=1}^m \mathbf{X}'_i \sim \text{Mult}(m, M, \mathbf{R}_D(h))$. Finally let $f : \Delta_M^m \rightarrow \mathbb{R}$ be defined by $f(x_1, \dots, x_m) := \exp\left(\beta d\left(\frac{1}{m} \sum_{i=1}^m x_i, \mathbf{R}_D(h)\right)\right)$, which is convex since the average is linear, d is convex and the exponential is non-decreasing and convex. Then, by swapping expectations (which is permitted by Fubini's theorem since the argument is non-negative) and applying Lemma 5, we have that $\mathbb{E}_{S' \sim D^m} Z_P(S')$ can be written as

$$\begin{aligned} \mathbb{E}_{S' \sim D^m} Z_P(S') &= \mathbb{E}_{S' \sim D^m} \mathbb{E}_{h \sim P} \exp\left(\beta d(\mathbf{R}_{S'}(h), \mathbf{R}_D(h))\right) \\ &= \mathbb{E}_{h \sim P} \mathbb{E}_{S' \sim D^m} \exp\left(\beta d(\mathbf{R}_{S'}(h), \mathbf{R}_D(h))\right) \\ &= \mathbb{E}_{h \sim P} \mathbb{E}_{\mathbf{X}_1, \dots, \mathbf{X}_m} \exp\left(\beta d\left(\frac{1}{m} \sum_{i=1}^m \mathbf{X}_i, \mathbf{R}_D(h)\right)\right) \\ &\leq \mathbb{E}_{h \sim P} \mathbb{E}_{\mathbf{X}'_1, \dots, \mathbf{X}'_m} \exp\left(\beta d\left(\frac{1}{m} \sum_{i=1}^m \mathbf{X}'_i, \mathbf{R}_D(h)\right)\right) \\ &= \mathbb{E}_{h \sim P} \mathbb{E}_{\mathbf{Y}} \exp\left(\beta d\left(\frac{1}{m} \mathbf{Y}, \mathbf{R}_D(h)\right)\right) \\ &= \mathbb{E}_{h \sim P} \sum_{\mathbf{k} \in S_{m, M}} \text{Mult}(\mathbf{k}; m, M, \mathbf{R}_D(h)) \exp\left(\beta d\left(\frac{\mathbf{k}}{m}, \mathbf{R}_D(h)\right)\right) \end{aligned}$$

$$\leq \sup_{\mathbf{r} \in \Delta_M} \left[\sum_{\mathbf{k} \in S_{m,M}} \text{Mult}(\mathbf{k}; m, M, \mathbf{r}) \exp\left(\beta d\left(\frac{\mathbf{k}}{m}, \mathbf{r}\right)\right) \right].$$

Which is the definition of $\mathcal{I}_d(m, \beta)$. Inequality (1) then follows by substituting this bound on $\mathbb{E}_{S' \sim D^m} Z_P(S')$ into (2) and dividing by β . \square

3.2 Statement and proof of the generalised corollary

We now apply our generalised theorem with $\beta = m$ and $d(\mathbf{q}, \mathbf{p}) = \text{kl}(\mathbf{q} \parallel \mathbf{p})$. This results in the following corollary, analogous to Corollary 2 (although the multi-dimensionality makes the proof much more involved, requiring multiple lemmas and extra arguments to make the main idea go through). We give two forms of the bound since, while the second is looser, the first is not practical to calculate except when m is very small.

Corollary 7. *Let \mathcal{X} and \mathcal{Y} be arbitrary sets and $\bigcup_{j=1}^M E_j$ be a disjoint partition of \mathcal{Y}^2 . Let $D \in \mathcal{M}(\mathcal{X} \times \mathcal{Y})$ be a data-generating distribution and $\mathcal{H} \subseteq \mathcal{M}(\mathcal{Y})^{\mathcal{X}}$ be a hypothesis class. For any prior $P \in \mathcal{M}(\mathcal{H})$, $\delta \in (0, 1]$ and sample size m , with probability at least $1 - \delta$ over the random draw $S \sim D^m$, we have that simultaneously for all posteriors $Q \in \mathcal{M}(\mathcal{H})$*

$$\text{kl}(\mathbf{R}_S(Q) \parallel \mathbf{R}_D(Q)) \leq \frac{1}{m} \left[\text{KL}(Q \parallel P) + \ln \left(\frac{m!}{\delta m^m} \sum_{\mathbf{k} \in S_{m,M}} \prod_{j=1}^M \frac{k_j^{k_j}}{k_j!} \right) \right] \quad (3)$$

$$\leq \frac{1}{m} \left[\text{KL}(Q \parallel P) + \ln \left(\frac{1}{\delta} \sqrt{\pi} e^{1/(12m)} \left(\frac{m}{2}\right)^{\frac{M-1}{2}} \sum_{z=0}^{M-1} \binom{M}{z} \frac{1}{(\pi m)^{z/2} \Gamma\left(\frac{M-z}{2}\right)} \right) \right], \quad (4)$$

where the second inequality holds provided $m \geq M$. Further, the bounds are unchanged if one restricts to an ordinary hypothesis class, namely if $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$.

While analogous corollaries can be obtained from Theorem 3 by other choices of convex function d , the kl-divergence leads to convenient cancellations that remove the dependence of $\mathcal{I}_{\text{kl}}(m, \beta, \mathbf{r})$ on \mathbf{r} , making $\mathcal{I}_{\text{kl}}(m, \beta) := \sup_{\mathbf{r} \in \Delta_M} \mathcal{I}_{\text{kl}}(m, \beta, \mathbf{r})$ simple to evaluate. Note (4) is logarithmic in $1/\delta$ (typical of PAC-Bayes bounds) and thus the confidence can be increased very cheaply. Ignoring logarithmic terms, (4) is $\mathcal{O}(1/m)$, also as expected. As for M , a simple analysis shows that (4) grows only sublinearly in M , meaning M can be made quite large provided one has a reasonable amount of data. To prove Corollary 7 we require Lemma 8, the proof of which is deferred to Appendix B.2.

Lemma 8. *For integers $M \geq 1$ and $m \geq M$, $\sum_{\mathbf{k} \in S_{m,M}^{>0}} \frac{1}{\prod_{j=1}^M \sqrt{k_j}} \leq \frac{\pi^{\frac{M}{2}} m^{\frac{M-2}{2}}}{\Gamma\left(\frac{M}{2}\right)}$.*

Proof. (of Corollary 7) Applying Theorem 3 with $d(\mathbf{q}, \mathbf{p}) = \text{kl}(\mathbf{q} \parallel \mathbf{p})$ (defined in Section 2) and $\beta = m$ gives that with probability at least $1 - \delta$ over $S \sim D^m$, simultaneously for all posteriors $Q \in \mathcal{M}(\mathcal{H})$, $\text{kl}(\mathbf{R}_S(Q) \parallel \mathbf{R}_D(Q)) \leq \frac{1}{m} [\text{KL}(Q \parallel P) + \ln \frac{\mathcal{I}_{\text{kl}}(m, m)}{\delta}]$, where $\mathcal{I}_{\text{kl}}(m, m) := \sup_{\mathbf{r} \in \Delta_M} [\sum_{\mathbf{k} \in S_{m,M}} \text{Mult}(\mathbf{k}; m, M, \mathbf{r}) \exp(m \text{kl}(\frac{\mathbf{k}}{m}, \mathbf{r}))]$. Thus, to establish the first inequality of the corollary, it suffices to show that

$$\mathcal{I}_{\text{kl}}(m, m) \leq \frac{m!}{m^m} \sum_{\mathbf{k} \in S_{m,M}} \prod_{j=1}^M \frac{k_j^{k_j}}{k_j!}. \quad (5)$$

To see this, for each fixed $\mathbf{r} = (r_1, \dots, r_M) \in \Delta_M$ let $J_{\mathbf{r}} = \{j \in [M] : r_j = 0\}$. Then $\text{Mult}(\mathbf{k}; m, M, \mathbf{r}) = 0$ for any $\mathbf{k} \in S_{m,M}$ such that $k_j \neq 0$ for some $j \in J_{\mathbf{r}}$. For the other $\mathbf{k} \in S_{m,M}$, namely those such that $k_j = 0$ for all $j \in J_{\mathbf{r}}$, the probability term can be written as $\text{Mult}(\mathbf{k}; m, M, \mathbf{r}) = \frac{m!}{\prod_{j=1}^M k_j!} \prod_{j=1}^M r_j^{k_j} = \frac{m!}{\prod_{j \notin J_{\mathbf{r}}} k_j!} \prod_{j \notin J_{\mathbf{r}}} r_j^{k_j}$, and (recalling the convention that $0 \ln \frac{0}{0} = 0$) the term $\exp(m \text{kl}(\frac{\mathbf{k}}{m}, \mathbf{r}))$ can be written as

$$\exp\left(m \sum_{j=1}^M \frac{k_j}{m} \ln \frac{k_j}{r_j}\right) = \exp\left(\sum_{j \notin J_{\mathbf{r}}} k_j \ln \frac{k_j}{m r_j}\right) = \prod_{j \notin J_{\mathbf{r}}} \left(\frac{k_j}{m r_j}\right)^{k_j} = \frac{1}{m^m} \prod_{j \notin J_{\mathbf{r}}} \left(\frac{k_j}{r_j}\right)^{k_j},$$

where the last equality is obtained by recalling that the k_j sum to m . Substituting these two expressions into the definition of $\mathcal{I}_{\text{kl}}(m, m)$ and only summing over those $\mathbf{k} \in S_{m, M}$ with non-zero probability, we obtain

$$\begin{aligned}
\sum_{\mathbf{k} \in S_{m, M}} \text{Mult}(\mathbf{k}; m, M, \mathbf{r}) \exp(m \text{kl}(\frac{\mathbf{k}}{m}, \mathbf{r})) &= \sum_{\substack{\mathbf{k} \in S_{m, M}: \\ \forall j \in J_{\mathbf{r}} k_j = 0}} \text{Mult}(\mathbf{k}; m, M, \mathbf{r}) \exp(m \text{kl}(\frac{\mathbf{k}}{m}, \mathbf{r})) \\
&= \sum_{\substack{\mathbf{k} \in S_{m, M}: \\ \forall j \in J_{\mathbf{r}} k_j = 0}} \frac{m!}{\prod_{j \notin J_{\mathbf{r}}} k_j!} \prod_{j \notin J_{\mathbf{r}}} r_j^{k_j} \frac{1}{m^m} \prod_{j \notin J_{\mathbf{r}}} \left(\frac{k_j}{r_j}\right)^{k_j} \\
&= \frac{m!}{m^m} \sum_{\substack{\mathbf{k} \in S_{m, M}: \\ \forall j \in J_{\mathbf{r}} k_j = 0}} \prod_{j \notin J_{\mathbf{r}}} \frac{k_j^{k_j}}{k_j!} \\
&= \frac{m!}{m^m} \sum_{\substack{\mathbf{k} \in S_{m, M}: \\ \forall j \in J_{\mathbf{r}} k_j = 0}} \prod_{j=1}^M \frac{k_j^{k_j}}{k_j!} \quad (\text{because } \frac{0^0}{0!} = 1) \\
&\leq \frac{m!}{m^m} \sum_{\mathbf{k} \in S_{m, M}} \prod_{j=1}^M \frac{k_j^{k_j}}{k_j!}.
\end{aligned}$$

Since this is independent of \mathbf{r} , it also holds after taking the supremum over $\mathbf{r} \in \Delta_M$ of the left hand side. We have thus established (5) and hence (3). Now, defining $f: \bigcup_{M=2}^{\infty} S_{m, M} \rightarrow \mathbb{R}$ by $f(\mathbf{k}) = \prod_{j=1}^{|\mathbf{k}|} k_j^{k_j} / k_j!$, we see that to establish inequality (4) it suffices to show that

$$\frac{m!}{m^m} \sum_{\mathbf{k} \in S_{m, M}} f(\mathbf{k}) \leq \sqrt{\pi} e^{1/12m} \left(\frac{m}{2}\right)^{\frac{M-1}{2}} \sum_{z=0}^{M-1} \binom{M}{z} \frac{1}{(\pi m)^{z/2} \Gamma\left(\frac{M-z}{2}\right)}. \quad (6)$$

We show this by upper bounding each $f(\mathbf{k})$ individually using Stirling's formula: $\forall n \geq 1 \sqrt{2\pi n} \left(\frac{n}{e}\right)^n < n! < \sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n}}$. Since we cannot use this to upper bound $1/k_j!$ when $k_j = 0$, we partition the sum above according to the number of coordinates of \mathbf{k} at which $k_j = 0$. Let z index the number of such coordinates. Since f is symmetric under permutations of its arguments,

$$\sum_{\mathbf{k} \in S_{m, M}} f(\mathbf{k}) = \sum_{z=0}^{M-1} \binom{M}{z} \sum_{\mathbf{k} \in S_{m, M-z}^{>0}} f(\mathbf{k}). \quad (7)$$

For $\mathbf{k} \in S_{m, M}^{>0}$ Stirling's formula yields $f(\mathbf{k}) \leq \prod_{j=1}^M \frac{k_j^{k_j}}{\sqrt{2\pi k_j} \left(\frac{k_j}{e}\right)^{k_j}} = \prod_{j=1}^M \frac{e^{k_j}}{\sqrt{2\pi k_j}} = \frac{e^m}{(2\pi)^{M/2}} \prod_{j=1}^M \frac{1}{\sqrt{k_j}}$. An application of Lemma 8 now gives

$$\sum_{\mathbf{k} \in S_{m, M-z}^{>0}} f(\mathbf{k}) \leq \frac{e^m}{(2\pi)^{M/2}} \sum_{\mathbf{k} \in S_{m, M-z}^{>0}} \prod_{j=1}^M \frac{1}{\sqrt{k_j}} \leq \frac{e^m}{(2\pi)^{M/2}} \frac{\pi^{\frac{M-z}{2}} m^{\frac{M-z-2}{2}}}{\Gamma\left(\frac{M-z}{2}\right)} = \frac{e^m m^{\frac{M-2}{2}}}{2^{\frac{M}{2}} (\pi m)^{z/2} \Gamma\left(\frac{M-z}{2}\right)}.$$

Substituting this into equation (7) and bounding $m!$ using Stirling's formula, we have

$$\begin{aligned}
\frac{m!}{m^m} \sum_{\mathbf{k} \in S_{m, M}} f(\mathbf{k}) &\leq \frac{\sqrt{2\pi m} e^{1/12m}}{e^m} \sum_{z=0}^{M-1} \binom{M}{z} \frac{e^m m^{\frac{M-2}{2}}}{2^{M/2} (\pi m)^{z/2} \Gamma\left(\frac{M-z}{2}\right)} \\
&= \sqrt{\pi} e^{1/12m} \left(\frac{m}{2}\right)^{\frac{M-1}{2}} \sum_{z=0}^{M-1} \binom{M}{z} \frac{1}{(\pi m)^{z/2} \Gamma\left(\frac{M-z}{2}\right)}
\end{aligned}$$

which is (6), establishing (4) and therefore completing the proof. \square

4 Implied Bounds and Construction of a Differentiable Training Objective

As already discussed, a multitude of bounds can be derived from Theorem 3 and Corollary 7, all of which then hold simultaneously with high probability. For example, suppose after a use of Corollary 7 we have a bound of the form $\text{kl}(\mathbf{R}_S(Q) \parallel \mathbf{R}_D(Q)) \leq B$. The following proposition then yields the bounds $L_j \leq R_D^j(Q) \leq U_j$, where $L_j := \inf\{p \in [0, 1] : \text{kl}(R_S^j(Q) \parallel p) \leq B\}$ and $U_j := \sup\{p \in [0, 1] : \text{kl}(R_S^j(Q) \parallel p) \leq B\}$. Moreover, since in the worst case we have $\text{kl}(\mathbf{R}_S(Q) \parallel \mathbf{R}_D(Q)) = B$, the proposition shows that the lower and upper bounds L_j and U_j are the tightest possible, since if $R_D^j(Q) \notin [L_j, U_j]$ then $\text{kl}(R_S^j(Q) \parallel R_D^j(Q)) > B$ implying $\text{kl}(\mathbf{R}_S(Q) \parallel \mathbf{R}_D(Q)) > B$. For a more precise version of this argument and a proof of Proposition 9, see Appendix B.3.

Proposition 9. *Let $\mathbf{q}, \mathbf{p} \in \Delta_M$. Then $\text{kl}(q_j \parallel p_j) \leq \text{kl}(\mathbf{q} \parallel \mathbf{p})$ for all $j \in [M]$, with equality when $p_i = \frac{1-p_j}{1-q_j} q_i$ for all $i \neq j$.*

As a second much more interesting example, suppose we can quantify how bad an error of each type is by means of a loss vector $\ell \in [0, \infty)^M$, where ℓ_j is the loss we attribute to an error of type E_j . We may then be interested in bounding the *total risk* $R_D^T(Q) \in [0, \infty)$ of Q which, recall, is defined by $R_D^T(Q) := \sum_{j=1}^M \ell_j R_D^j(Q)$. Indeed, given a bound of the form $\text{kl}(\mathbf{R}_S(Q) \parallel \mathbf{R}_D(Q)) \leq B$, we can derive $R_D^T(Q) \leq \sup\{\sum_{j=1}^M \ell_j r_j : \mathbf{r} \in \Delta_M, \text{kl}(\mathbf{R}_S(Q) \parallel \mathbf{r}) \leq B\}$. This motivates the following definition of $\text{kl}_\ell^{-1}(\mathbf{u} \mid c)$. To see that this is indeed well-defined (at least when $\mathbf{u} \in \Delta_M^{\geq 0}$), see the discussion at the beginning of Appendix B.4.

Definition 10. *For $\mathbf{u} \in \Delta_M, c \in [0, \infty)$ and $\ell \in [0, \infty)^M$, define $\text{kl}_\ell^{-1}(\mathbf{u} \mid c) = \sup\{\sum_{j=1}^M \ell_j v_j : \mathbf{v} \in \Delta_M, \text{kl}(\mathbf{u} \parallel \mathbf{v}) \leq c\}$.*

Can we calculate $\text{kl}_\ell^{-1}(\mathbf{u} \mid c)$ and hence $f_\ell(\text{kl}_\ell^{-1}(\mathbf{u} \mid c))$ in order to evaluate the bound on the total risk? Additionally, if we wish to use the bound on the total risk as a training objective, can we calculate the partial derivatives of $f_\ell^*(\mathbf{u}, c) := f_\ell(\text{kl}_\ell^{-1}(\mathbf{u} \mid c))$ with respect to the u_j and c so that we can use gradient descent? Our Proposition 11 answers both of these questions in the affirmative, at least in the sense that it provides a speedy method for approximating these quantities to arbitrary precision provided $u_j > 0$ for all $j \in [M]$ and $c > 0$. Indeed, the only approximation step required is that of approximating the unique root of a continuous and strictly increasing scalar function. Thus, provided the u_j themselves are differentiable, Corollary 7 combined with Proposition 11 yields a tractable and fully differentiable objective that can be used for training. More details on how this can be done, including an algorithm written in pseudocode, can be found in Appendix A. While somewhat analogous to the technique used in Clerico et al. (2022) to obtain derivatives of the one-dimensional kl-inverse, our proposition directly yields derivatives on the total risk by (implicitly) employing the envelope theorem (see for example Takayama & Akira, 1985). Since the proof of Proposition 11 is rather long and technical, we defer it to Appendix B.4.

Proposition 11. *Fix $\ell \in [0, \infty)^M$ such that not all ℓ_j are equal, and define $f_\ell : \Delta_M \rightarrow [0, \infty)$ by $f_\ell(\mathbf{v}) := \sum_{j=1}^M \ell_j v_j$. For all $\tilde{\mathbf{u}} = (\mathbf{u}, c) \in \Delta_M^{\geq 0} \times (0, \infty)$, define $\mathbf{v}^*(\tilde{\mathbf{u}}) := \text{kl}_\ell^{-1}(\mathbf{u} \mid c) \in \Delta_M$ and let $\mu^*(\tilde{\mathbf{u}}) \in (-\infty, -\max_j \ell_j)$ be the unique solution to $c = \phi_\ell(\mu)$, where $\phi_\ell : (-\infty, -\max_j \ell_j) \rightarrow \mathbb{R}$ is given by $\phi_\ell(\mu) := \ln(-\sum_{j=1}^M \frac{u_j}{\mu + \ell_j}) + \sum_{j=1}^M u_j \ln(-(\mu + \ell_j))$, which is continuous and strictly increasing. Then $\mathbf{v}^*(\tilde{\mathbf{u}}) = \text{kl}_\ell^{-1}(\mathbf{u} \mid c)$ is given by*

$$\mathbf{v}^*(\tilde{\mathbf{u}})_j = \frac{\lambda^*(\tilde{\mathbf{u}}) u_j}{\mu^*(\tilde{\mathbf{u}}) + \ell_j} \quad \text{for } j \in [M], \quad \text{where } \lambda^*(\tilde{\mathbf{u}}) = \left(\sum_{j=1}^M \frac{u_j}{\mu^*(\tilde{\mathbf{u}}) + \ell_j} \right)^{-1}.$$

Further, defining $f_\ell^* : \Delta_M^{\geq 0} \times (0, \infty) \rightarrow [0, \infty)$ by $f_\ell^*(\tilde{\mathbf{u}}) := f_\ell(\mathbf{v}^*(\tilde{\mathbf{u}}))$, we have that

$$\frac{\partial f_\ell^*}{\partial u_j}(\tilde{\mathbf{u}}) = \lambda^*(\tilde{\mathbf{u}}) \left(1 + \ln \frac{u_j}{\mathbf{v}^*(\tilde{\mathbf{u}})_j} \right) \quad \text{and} \quad \frac{\partial f_\ell^*}{\partial c}(\tilde{\mathbf{u}}) = -\lambda^*(\tilde{\mathbf{u}}).$$

5 Perspectives

By abstracting to a general setting of discretised error types, we established a novel type of generalisation bound (Theorem 3) providing far richer information than existing PAC-Bayes bounds. Through our Corollary

7 and Proposition 11, our bound inspires a training algorithm (see Appendix A) suitable for many different learning problems, including structured output prediction (as investigated by Cantelobre et al., 2020, in the PAC-Bayes setting), multi-task learning and learning-to-learn (see e.g. Maurer et al., 2016). We will demonstrate these applications and our bound’s utility for real-world learning problems in an empirical follow-up study. Note we require i.i.d. data, which in practice is frequently not the case or is hard to verify. Further, the number of error types M must be finite. While in continuous scenarios it would be preferable to be able to quantify the entire distribution of loss values without having to discretise into finitely many error types, in the multiclass setting our framework is entirely suitable.

References

- Pierre Alquier. User-friendly introduction to PAC-Bayes bounds. *arXiv preprint arXiv:2110.11216*, 2021.
- Amiran Ambroladze, Emilio Parrado-Hernández, and John Shawe-Taylor. Tighter PAC-Bayes bounds. In Bernhard Schölkopf, John C. Platt, and Thomas Hofmann (eds.), *Advances in Neural Information Processing Systems 19, Proceedings of the Twentieth Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 4-7, 2006*, pp. 9–16. MIT Press, 2006. URL <https://proceedings.neurips.cc/paper/2006/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>.
- Loubna Benabbou and Pascal Lang. PAC-Bayesian generalization bound for multi-class learning. In *NIPS 2017 Workshop. (Almost) 50 Shades of Bayesian Learning: PAC-Bayesian trends and insights*, 2017. URL https://bguedj.github.io/nips2017/pdf/PAC-Bayes_2017_paper_3.pdf.
- Felix Biggs and Benjamin Guedj. Differentiable PAC-Bayes objectives with partially aggregated neural networks. *Entropy*, 23(10):1280, 2021. doi: 10.3390/e23101280. URL <https://doi.org/10.3390/e23101280>.
- Felix Biggs and Benjamin Guedj. On margins and derandomisation in PAC-Bayes. In *AISTATS*, 2022a. URL <https://arxiv.org/abs/2107.03955>.
- Felix Biggs and Benjamin Guedj. Non-vacuous generalisation bounds for shallow neural networks. *arXiv preprint arXiv:2202.01627*, 2022b.
- Luc Bégin, Pascal Germain, François Laviolette, and Jean-François Roy. PAC-Bayesian Bounds based on the Rényi Divergence. In Arthur Gretton and Christian C. Robert (eds.), *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pp. 435–444, Cadiz, Spain, 09–11 May 2016. PMLR. URL <https://proceedings.mlr.press/v51/begin16.html>.
- Théophile Cantelobre, Benjamin Guedj, María Pérez-Ortiz, and John Shawe-Taylor. A pac-bayesian perspective on structured prediction with implicit loss embeddings. *arXiv preprint arXiv:2012.03780*, 2020.
- Olivier Catoni. A PAC-Bayesian approach to adaptive classification. *preprint*, 840, 2003.
- Olivier Catoni. *Statistical Learning Theory and Stochastic Optimization: Ecole d’Eté de Probabilités de Saint-Flour XXXI-2001*. Springer, 2004.
- Olivier Catoni. *PAC-Bayesian Supervised Classification: The Thermodynamics of Statistical Learning*, volume 56 of *Institute of Mathematical Statistics (IMS) Lecture Notes - Monograph Series*. Institute of Mathematical Statistics, 2007. ISBN 9780940600720. URL <https://books.google.fr/books?id=acnaAAAAMAAJ>.
- Eugenio Clerico, George Deligiannidis, and Arnaud Doucet. Conditionally gaussian pac-bayes. In *International Conference on Artificial Intelligence and Statistics*, pp. 2311–2329. PMLR, 2022.
- Imre Csiszár. I-divergence geometry of probability distributions and minimization problems. *The Annals of Probability*, pp. 146–158, 1975.
- MD Donsker and SRS Varadhan. Large deviations for Markov processes and the asymptotic evaluation of certain markov process expectations for large times. In *Probabilistic Methods in Differential Equations*, pp. 82–88. Springer, 1975.

- Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. In *Conference on Uncertainty in Artificial Intelligence [UAI]*, 2017.
- Gintare Karolina Dziugaite and Daniel M. Roy. Entropy-SGD optimizes the prior of a PAC-Bayes bound: Generalization properties of entropy-SGD and data-dependent priors. In Jennifer G. Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1376–1385. PMLR, 2018. URL <http://proceedings.mlr.press/v80/dziugaite18a.html>.
- Gintare Karolina Dziugaite, Kyle Hsu, Waseem Gharbieh, Gabriel Arpino, and Daniel Roy. On the role of data in PAC-Bayes. In Arindam Banerjee and Kenji Fukumizu (eds.), *The 24th International Conference on Artificial Intelligence and Statistics, AISTATS 2021, April 13-15, 2021, Virtual Event*, volume 130 of *Proceedings of Machine Learning Research*, pp. 604–612. PMLR, 2021. URL <http://proceedings.mlr.press/v130/karolina-dziugaite21a.html>.
- Vasilii Feofanov, Emilie Devijver, and Massih-Reza Amini. Transductive bounds for the multi-class majority vote classifier. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33:3566–3573, 2019. doi: 10.1609/aaai.v33i01.33013566. URL <https://ojs.aaai.org/index.php/AAAI/article/view/4236>.
- Benjamin Guedj. A Primer on PAC-Bayesian Learning. In *Proceedings of the second congress of the French Mathematical Society*, 2019. URL <https://arxiv.org/abs/1901.05353>.
- Sokol Koço and Cécile Capponi. On multi-class classification through the minimization of the confusion matrix norm. In Cheng Soon Ong and Tu Bao Ho (eds.), *Proceedings of the 5th Asian Conference on Machine Learning*, volume 29 of *Proceedings of Machine Learning Research*, pp. 277–292, Australian National University, Canberra, Australia, 13–15 Nov 2013. PMLR. URL <https://proceedings.mlr.press/v29/Koco13.html>.
- François Laviolette, Emilie Morvant, Liva Ralaivola, and Jean-François Roy. Risk upper bounds for general ensemble methods with an application to multiclass classification. *Neurocomputing*, 219:15–25, 2017. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2016.09.016>. URL <https://www.sciencedirect.com/science/article/pii/S0925231216310177>.
- Gaël Letarte, Pascal Germain, Benjamin Guedj, and François Laviolette. Dichotomize and generalize: PAC-Bayesian binary activated deep neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. dAlché Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp. 6872–6882. Curran Associates, Inc., 2019.
- Guy Lever, François Laviolette, and John Shawe-Taylor. Distribution-dependent PAC-Bayes priors. In *International Conference on Algorithmic Learning Theory*, pp. 119–133. Springer, 2010.
- Guy Lever, François Laviolette, and John Shawe-Taylor. Tighter PAC-Bayes bounds through distribution-dependent priors. *Theoretical Computer Science*, 473:4–28, February 2013. ISSN 0304-3975. doi: 10.1016/j.tcs.2012.10.013. URL <https://linkinghub.elsevier.com/retrieve/pii/S0304397512009346>.
- Andreas Maurer. A note on the PAC-Bayesian theorem. *arXiv preprint cs/0411099*, 2004.
- Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning. *J. Mach. Learn. Res.*, 17:81:1–81:32, 2016. URL <http://jmlr.org/papers/v17/15-242.html>.
- David A McAllester. Some PAC-Bayesian theorems. In *Proceedings of the eleventh annual conference on Computational Learning Theory*, pp. 230–234. ACM, 1998.
- David A McAllester. PAC-Bayesian model averaging. In *Proceedings of the twelfth annual conference on Computational Learning Theory*, pp. 164–170. ACM, 1999.

- Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte carlo gradient estimation in machine learning. *J. Mach. Learn. Res.*, 21(132):1–62, 2020.
- Emilie Morvant, Sokol Koço, and Liva Ralaivola. PAC-Bayesian generalization bound on confusion matrix for multi-class classification. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*. icml.cc / Omnipress, 2012. URL <http://icml.cc/2012/papers/434.pdf>.
- Behnam Neyshabur, Srinadh Bhojanapalli, and Nathan Srebro. A PAC-Bayesian approach to spectrally-normalized margin bounds for neural networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL https://openreview.net/forum?id=Skz_WfbCZ.
- Emilio Parrado-Hernández, Amiran Ambroladze, John Shawe-Taylor, and Shiliang Sun. PAC-Bayes bounds with data dependent priors. *J. Mach. Learn. Res.*, 13:3507–3531, 2012. URL <http://dl.acm.org/citation.cfm?id=2503353>.
- María Pérez-Ortiz, Omar Rivasplata, Benjamin Guedj, Matthew Gleeson, Jingyu Zhang, John Shawe-Taylor, Mirosław Bober, and Josef Kittler. Learning pac-bayes priors for probabilistic neural networks. *arXiv preprint arXiv:2109.10304*, 2021.
- Maria Perez-Ortiz, Omar Rivasplata, John Shawe-Taylor, and Csaba Szepesvari. Tighter risk certificates for neural networks. *Journal of Machine Learning Research*, 22(227):1–40, 2021. URL <http://jmlr.org/papers/v22/20-879.html>.
- Omar Rivasplata, Csaba Szepesvári, John Shawe-Taylor, Emilio Parrado-Hernández, and Shiliang Sun. PAC-Bayes bounds for stable algorithms with instance-dependent priors. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 9234–9244, 2018. URL <https://proceedings.neurips.cc/paper/2018/hash/386854131f58a556343e056f03626e00-Abstract.html>.
- Matthias Seeger. PAC-Bayesian generalisation error bounds for Gaussian process classification. *Journal of Machine Learning Research*, 3(Oct):233–269, 2002.
- Akira Takayama and Takayama Akira. *Mathematical economics*. Cambridge university press, 1985.
- 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 *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=BJgqqsAct7>.