

# Safe Learning: *bridging the gap between Bayes, MDL and statistical learning theory via empirical convexity*

Peter Grünwald

PDG@CWI.NL

CWI Amsterdam and Mathematical Institute, Leiden University, The Netherlands

**Editor:** Sham Kakade, Ulrike von Luxburg

## Abstract

We extend Bayesian MAP and Minimum Description Length (MDL) learning by testing whether the data can be substantially more compressed by a mixture of the MDL/MAP distribution with another element of the model, and adjusting the learning rate if this is the case. While standard Bayes and MDL can fail to converge if the model is wrong, the resulting “safe” estimator continues to achieve good rates with wrong models. Moreover, when applied to classification and regression models as considered in statistical learning theory, the approach achieves optimal rates under, e.g., Tsybakov’s conditions, and reveals new situations in which we can penalize by  $(-\log \text{PRIOR})/n$  rather than  $\sqrt{(-\log \text{PRIOR})/n}$ .

## 1. Introduction

*1. Learning Theory; Predictor Models.* In much of statistical learning and machine learning theory, the goal is to learn, based on a set of observed data  $Z^n = (Z_1, \dots, Z_n)$ , a predictor  $\check{f}$  taken from some set of candidate prediction rules  $\mathcal{F}$ . Here each  $Z_i = (X_i, Y_i)$ , each  $X_i$  takes values in some set  $\mathcal{X}$ , each  $Y_i$  takes values in  $\mathcal{Y}$ , and  $\mathcal{F}$  is a set of functions  $f: \mathcal{X} \rightarrow \mathcal{Y}$ . The  $Z_i$  are assumed to be sampled i.i.d. according to some distribution  $P^*$  on  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ . The learned predictor  $\check{f}$  should have a small *generalization error* or *risk*, defined as  $\text{RISK}(f) := E^*[\text{LOSS}(Y, f(X))]$  where  $\text{LOSS}$  is some given loss function and here, as elsewhere in this paper,  $E^* = E_{(X,Y) \sim P^*}$  denotes joint expectation of  $(X, Y)$  over  $P^*$ . In a typical classification setting,  $\mathcal{Y} = \{0, 1\}$  and  $\text{LOSS}(y, \hat{y}) := |y - \hat{y}|$  is the 0/1-loss; in typical regression problems,  $\mathcal{Y} = \mathbb{R}$  and  $\text{LOSS}(y, \hat{y}) := (y - \hat{y})^2$  is the squared loss. Crucially, risk bounds are usually proved in worst-case settings, using only weak assumptions on  $P^*$ .

*2. Standard Statistics; Probability Models.* Here one models uncertainty by a statistical model, i.e. a set of probability distributions  $\mathcal{P}$ , and the goal is to learn a distribution  $\check{p}$  that is a good representation of the underlying distribution  $P^*$  from which the data  $Z^n$  are sampled. Here we focus on the case that  $Z^n$  are i.i.d.,  $Z_i = (X_i, Y_i)$  as above, and  $\mathcal{P}$  is a set of conditional distributions  $p(y | x)$ , identified by their mass functions (if  $\mathcal{Y}$  is finite/countable) or densities, and extended to  $n$  outcomes by independence. Witness papers such as *The Two Cultures* (Breiman, 2001), the difference between statistical/machine learning theory and standard statistics based on probability models is often regarded as fundamental. Here, I propose a first, preliminary, attempt at an overarching, single theory of learning, as embodied by a new ‘safe’ estimator. It is called ‘safe’ because, when applied to probability models  $\mathcal{P}$ , then, unlike standard Bayes and MDL, it is guaranteed to perform well in the often inevitable situation that ‘all models (elements of  $\mathcal{P}$ ) are wrong, yet some are useful’.

**Safe Estimation for Probability Models** For probability models  $\mathcal{P}$ , the safe estimator behaves similarly to the Bayesian MAP or two-part MDL estimator. Following Barron and Cover (1991), we define the  $\kappa$ -two part estimator, written as  $\check{p}_\kappa$ , as a generalization of the MAP/MDL estimator, as follows: fix some prior distribution  $w$  and some  $\kappa > 0$ . For each  $x^n \in \mathcal{X}^n, y^n \in \mathcal{Y}^n$ ,  $\check{p}_\kappa$  is defined\* as the  $p \in \mathcal{P}$  achieving<sup>1</sup>

$$\min_{p \in \mathcal{P}} \{ -\kappa \log w(p) - \log p(y^n | x^n) \}. \quad (1)$$

When  $\kappa \geq 1$ , then, via the Kraft inequality, (1) can be thought of as the number of bits needed to encode  $Y^n$  given  $X^n$  in a two-stage code;  $-\kappa \log w(p)$  is the codelength needed to encode  $p$ , and acts as a complexity penalty.  $-\log p(y^n | x^n)$  is the codelength of the data  $y^n$  when encoded with the help of  $p$  and  $x^n$ . To get good convergence rates, one needs to set  $\kappa > 1$  (Zhang, 2006); while any fixed  $\kappa > 1$  will do, for ‘standard 2-part MDL’ one takes  $\kappa = 2$  which is mathematically convenient (Barron and Cover, 1991). In contrast, the safe estimator is defined (in Section 2, Eq. (9)) as  $\check{p}_{\text{SAFE}} = \check{p}_{2\check{\kappa}_{\text{SAFE}}}$  where  $\check{\kappa}_{\text{SAFE}}$  is not fixed but determined by the data.  $\check{\kappa}_{\text{SAFE}}$  will be a small constant  $\geq 2$ , *unless* the data indicate that the model is misspecified (wrong). Whereas ordinary Bayesian and MDL approaches can fail to converge in this case (Example 7 below), the safe estimator continues to perform well in the following sense: suppose that data  $Z^n$  are i.i.d.  $\sim P^*$ , as above, where for each  $x$ ,  $P^*(Y = \cdot | X = x)$  has conditional density  $p^*(\cdot | x)$ . Let  $q$  be the best approximation within  $\mathcal{P}$  of  $p^*$  in terms of Kullback-Leibler (KL) divergence. Then the KL divergence between  $\check{p}_{\text{SAFE}}$  and  $p^*$  converges to the KL divergence between  $q$  and  $p^*$  at fast rates. To express this formally, for any two conditional densities  $p$  and  $p'$ , we define\* the *generalized KL divergence* (Grünwald, 2007) relative to  $P^*$  as

$$D^*(p' \| p) := E^*[-\log p(Y | X) + \log p'(Y | X)].$$

Then, for  $q$  satisfying  $\inf_{p \in \mathcal{P}} D^*(p^* \| p) = D^*(p^* \| q)$  we prove, under suitable regularity conditions, in Theorem 1 in combination with Theorem 3 below that  $D(p^* \| \check{p}_{\text{SAFE}}) - D(p^* \| q) \rightarrow 0$ , or equivalently,  $D^*(q \| \check{p}_{\text{SAFE}}) \rightarrow 0$ , in probability at fast rates.

**Safe Estimation for Predictor Models** In our overarching approach, all models are formally defined as sets of probability distributions. Predictor models  $\mathcal{F}$  are “transformed” into corresponding probability models  $\mathcal{P}_{\mathcal{F}} := \{p_f | f \in \mathcal{F}\}$  by a standard transformation (called ‘entropification’ and extensively motivated from an MDL perspective by Grünwald (1999)): for each  $f \in \mathcal{F}$ ,

$$p_f(y | x) := \frac{1}{Z(\beta)} e^{-\beta \text{LOSS}(y, f(x))} \quad , \quad p_f(y^n | x^n) := \prod_{i=1}^n p_f(y_i | x_i). \quad (2)$$

Here  $Z(\beta) = \int_{y \in \mathcal{Y}} \exp(\beta \text{LOSS}(y, f(x))) dy$  is a normalization factor (if  $\mathcal{Y}$  is finite/countable, then here, as everywhere else in this paper, the integral should be replaced by a sum). In this preliminary study, we set  $\beta$  to some fixed value, say, 1 (but see Section 6). For the squared loss,  $Z(\beta)$  does not depend on  $f(x)$ ; if we set  $\beta = 1/2\sigma^2$ , we see that (2) expresses that  $Y$  is Gaussian with mean  $f(X)$  and variance  $\sigma^2$ . For the 0/1-loss,  $Z(\beta)$  does not

1. Distracting aspects of proofs (such as showing that the minimum of a function exists) have been omitted in this paper, but will be provided in the journal version. Such details are marked by a \*, such as here\*.

depend on  $f(x)$  either; loss functions for which  $Z(\beta)$  depends on  $f$  are handled as described under Eq.(7) below. Taking logarithms in (2), we then get that the *excess risk* of any  $f$  as compared to any  $g$  is a linear function of the generalized KL divergence of the corresponding distributions:

$$\begin{aligned} \text{EXCESS-RISK}(g\|f) &= \text{RISK}(f) - \text{RISK}(g) = \\ &E^*[\text{LOSS}(Y, f(X)) - \text{LOSS}(Y, g(X))] = \frac{1}{\beta}D^*(p_g\|p_f). \end{aligned} \quad (3)$$

Now let  $g$  be such that  $\text{RISK}(g) = \inf_{f \in \mathcal{F}} \text{RISK}(f)$ . Even if  $\mathcal{F}$  is a good ‘model’, i.e.  $\text{RISK}(g)$  is small, the corresponding model  $\mathcal{P}_{\mathcal{F}}$  will typically be misspecified (e.g. in the squared loss case, the ‘true’ noise may not be Gaussian at all). Since the safe estimator is immune to this problem, we can safely apply it to the model  $\mathcal{P}_{\mathcal{F}}$ . Then Theorems 1 and 3 show that the excess risk  $\text{EXCESS-RISK}(g\|\check{f}_{\text{SAFE}})$  converges to 0 at rates that are in many cases optimal; here  $\check{f}_{\text{SAFE}} := f$  for the  $f$  with  $\check{p}_{\text{SAFE}} = p_f$ . Thus, by the construction (2), convergence in generalized KL-divergence becomes equivalent to convergence in the loss function of interest.

**The Role of Convexity** Our starting point is the known fact that ‘standard MDL still works’, i.e., (broadly speaking),  $D^*(q\|\check{p}_2) \rightarrow 0$  at the appropriate rate if the closure  $\langle \mathcal{P} \rangle$  (suitably defined as in (10) below) of the model  $\mathcal{P}$  is *convex* (Li, 1999, Theorem 5.5); see also (Kleijn and van der Vaart, 2006). Our first observation is that, even if  $\langle \mathcal{P} \rangle$  is not convex, then as long as we have the weaker condition

$$\inf_{p \in \mathcal{P}} D^*(p^*\|p) = \inf_{p \in \text{CONVEX-HULL}(\mathcal{P})} D^*(p^*\|p), \quad (4)$$

we still get that  $D^*(q\|\check{p}_2) \rightarrow 0$  at the right rates. Now define, for  $\eta \leq 1$ , the model  $\mathcal{P}^{(\eta)} := \{p^{(\eta)} \mid p \in \mathcal{P}\}$ , where  $p^{(\eta)}(y \mid x) \propto (p(y \mid x))^\eta$  (for predictor models  $\mathcal{P}_{\mathcal{F}}$ , this corresponds to replacing  $\beta$  in (2) by  $\eta \cdot \beta$ ; a precise definition is beneath (7) below). Our second insight is that, even if (4) does not hold for  $\mathcal{P}$ , then still, for all  $\eta$  no greater than some critical value  $\eta_{\text{CRIT}}$ , (4) will actually hold with  $\mathcal{P}$  replaced by  $\mathcal{P}^{(\eta)}$  and thus ‘standard MDL still works’ for  $\mathcal{P}^{(\eta)}$ . The third insight is that the MDL estimator  $\check{p}_{2\kappa}$  for model  $\mathcal{P}$  with  $\kappa = 1/\eta_{\text{CRIT}}$  is essentially equivalent to the standard MDL estimator  $\check{p}_2^{(\eta)}$  for model  $\mathcal{P}^{(\eta_{\text{CRIT}})}$ ; indeed, we will prove (implicitly in Theorem 3) that the MDL estimator  $\check{p}_{2\kappa}$  with  $\kappa = 1/\eta_{\text{CRIT}}$  leads to good results for the model  $\mathcal{P}$ . The fourth, and main, insight is that, for any given  $\eta$ , we can *test* whether  $\eta \leq \eta_{\text{CRIT}}$ , i.e. whether (4) is the case for  $\mathcal{P}^{(\eta)}$ , by looking at the observed data: *essentially, the likelihood of the data according to  $\check{p}_2^{(\eta)}$  will be significantly smaller than the likelihood according to a 2-component mixture of  $\check{p}_2^{(\eta)}$  and another, suitably chosen  $p \in \mathcal{P}$ , if and only if (4) does not hold.* The minus logarithm of this discrepancy is measured in terms of a function CONV-LACK, a key concept of this paper, defined formally in (8). The safe estimator is defined as  $\check{p}_{\text{SAFE}} = \check{p}_{2\check{\kappa}_{\text{SAFE}}}$ , i.e. it is the  $2\check{\kappa}_{\text{SAFE}}$ -two-part estimator where  $\check{\kappa}_{\text{SAFE}} = 1/\eta$ , and  $\eta$  is determined by the data: it is effectively set to the largest value for which CONV-LACK is small, i.e. for which we cannot fit the data better by a two-component mixture.

**Overview of Results** In Section 2 we formally define CONV-LACK and  $\check{p}_{\text{SAFE}}$ . In Section 3 and 4,  $\check{p}_{\text{SAFE}}$  will be shown to converge at optimal rates up to log factors in a variety of

situations, as illustrated by examples in Section 4. Convergence of  $\check{p}_{\text{SAFE}}$  is shown in two stages: Theorem 1 bounds  $D^*(q\|\check{p})$  for arbitrary estimators  $\check{p}$  in terms of a ‘redundancy’ term RED and the CONV-LACK term. The redundancy term also shows up in classical MDL analyses and tends to 0 if we set  $\check{p}$  to a two-part estimator. Theorem 3 shows (essentially) that if  $\check{p}$  is set to  $\check{p}_{2\kappa}$  with  $\kappa \geq 1/\eta_{\text{CRIT}}$ , then the CONV-LACK-term is small as well. Taken together, Theorem 1 and 3 imply that  $\check{p}_{\text{SAFE}}$  converges (a) at the right rates if the model is correct or convex; and (b) also if the model is “incorrect in the worst possible manner”, and finally, (c) also if the model is a classification model, i.e.  $\mathcal{P} = \mathcal{P}_{\mathcal{F}}$  for some set of classifiers  $\mathcal{F}$ , and a Tsybakov condition holds for  $\mathcal{F}$ .

In case (a),  $\eta_{\text{CRIT}} = 1$ . As shown in Example 4, if the model is in fact correct ( $p^* = q$ ) we get the same bound on  $D^*(p^*\|\check{p}_{\text{SAFE}})$  as the bounds obtained on  $D^*(p^*\|\check{p}_2)$  by Barron and Cover (1991), but with a larger constant factor — this is the price we have to pay for using a method that still works if the model is incorrect in a situation in which the model in fact, turns out to be correct. In the special case that  $w(p^*) = w(q) > 0$ ,  $D^*(q\|\check{p}_{\text{SAFE}})$  will tend to 0 as  $O((\log n)/n)$ . In case (b),  $\eta_{\text{CRIT}}$  may be as small as  $C/\sqrt{n}$  for some  $C > 0$ . Example 6 shows that  $D^*(q\|\check{p}_{\text{SAFE}})$  may then tend to 0 at rate as slow as  $(\log n)/\sqrt{n}$ , a worst-case bound familiar from the statistical learning literature. In case (c),  $\eta_{\text{CRIT}} \asymp n^{-(1-\nu)/(2-\nu)}$  for some  $0 \leq \nu \leq 1$  and the convergence rate depends on  $\nu$ . In the special case that  $w(q) > 0$ ,  $D^*(q\|\check{p}_{\text{SAFE}})$  will tend to 0 as  $O((\log n)n^{-1/(2-\nu)})$ ; see above Example 7. The examples illustrate the generality of  $\check{p}_{\text{SAFE}}$ , capturing both the common asymptotics for density estimation if the model is correct and for statistical classifier learning under the celebrated Tsybakov condition.  $\check{p}_{\text{SAFE}}$  also gives a new interpretation of the difference in complexity penalties prescribed by MDL/Bayes on the one hand and learning theory approaches (such as Structural Risk Minimization and PAC-Bayesian methods) on the other: since the  $\mathcal{P}_{\mathcal{F}}$  constructed from a predictor model  $\mathcal{F}$  is in general nonconvex, we may have  $\eta_{\text{CRIT}} \ll 1$ , and the standard MDL penalties become too small.

A third main result is Theorem 2, which gives a new PAC-Bayes style empirical generalization bound in which, if we are ‘lucky’ on the observed data, the codelength ( $-\log w(\check{p}_{\text{SAFE}})$ ) only appears in an  $O(1/n)$  rather than an  $O(1/\sqrt{n})$  term, even if the empirical error of the learned classifier is not close to 0. As such it provides another step in the race to “root out the square root” that characterizes so much of the work on classification bounds in learning theory.

**Related Work — Learning the Learning Rate** The larger  $\eta = 1/\kappa$ , the more influence the data has on the chosen hypothesis  $\check{p}_{2\kappa}$ . For predictor models, the same holds for the  $\beta$  appearing in (2). Thus  $\eta \cdot \beta$  may be viewed as the learning rate. A straightforward approach to learn it from data is to fix  $\eta$  and instead pick the  $\beta$  in (2) that minimizes overall description length of the data, as suggested by Grünwald (1999). Soon after publishing that paper, it became clear to me that this does not work (this was shown formally in (Grünwald and Langford, 2007)), and I started looking for an estimator that performs as well as if the optimal learning rate  $\eta_{\text{CRIT}}$  had been known in advance. The safe estimator does achieve this goal, thus ending a twelve-year long search. In a sense,  $\check{p}_{\text{SAFE}}$  learns the optimal learning rate. Note however that we cannot prove that the  $\check{\kappa}_{\text{SAFE}}$  selected by the safe estimator is equal or close to  $1/\eta_{\text{CRIT}}$ ; we can only prove that it leads to the same asymptotic performance bounds.

A transformation similar to (2) is done in PAC-Bayesian methods (McAllester, 2003), where Bayesian averages of  $p_f$  are viewed as Gibbs distributions. Our approach is similar, yet closer in spirit to standard Bayes and MDL — There may be some relation with the advanced PAC-Bayesian analyses of Audibert (2004); Catoni (2007), who provide algorithms for predictor models that learn a learning rate (similar to  $\eta = 1/\kappa$ ) determined by the amount of ‘disagreement’ on the input data  $(X_1, \dots, X_n)$  between the chosen predictor  $\check{f}$  and other predictors in  $\mathcal{F}$ . It would be interesting to study the connections further.

Finally, our approach can (broadly) be seen as equipping (a form of) Bayesian inference with a frequentist test, and adjusting the priors if the test indicates that the model is misspecified. Such an idea was already suggested in broad terms by other researchers, e.g. Dawid (1982). It can also be viewed as equipping (a form of) MDL with a randomness test (can we compress the data more by stepping outside the model?), an idea that goes back to the Kolmogorov complexity roots of MDL.

## 2. The Safe Estimator

**Preliminaries** A *probabilistic model*  $\mathcal{P}$  is a *countable* set of conditional distributions on  $\mathcal{Y}$  given  $\mathcal{X}$ , identified with their mass functions (in case  $\mathcal{Y}$  is finite or countable) or otherwise their densities relative to Lebesgue measure, which we assume to exist;  $\mathcal{X}$  can be arbitrary. We allow the  $p \in \mathcal{P}$  to be defective (sum to less than one). That is,  $p(y | x)$  can be any function such that, for all  $x$ , for all  $y$ ,  $p(y | x) \geq 0$  and  $\sum_{y \in \mathcal{Y}} p(y | x) \leq 1$ . We extend  $p$  to  $n$  outcomes by independence, i.e.  $p(y^n | x^n) := \prod_{i=1}^n p(y_i | x_i)$ . For given  $z^n = (x_1, y_1), \dots, (x_n, y_n)$ ,  $z_i^n$  is shorthand for  $y^n | x^n$ , i.e.  $p(z_i^n) = p(y^n | x^n)$ . All logarithms are to base  $e$ .

Crucially, the models  $\mathcal{P}$  we are to consider, though countable, will usually represent very “large”, “complex” sets of distributions, which may be thought of as dense (in the information closure sense, see Section 3) subsets of an even larger, “nonparametric”  $\bar{\mathcal{P}}$  with  $\mathcal{P} \subset \bar{\mathcal{P}}$ : for example, we may consider the set  $\bar{\mathcal{P}}$  of all Gaussian mixtures with an arbitrary number of components, and then define  $\mathcal{P}$  as the subset of all  $p \in \bar{\mathcal{P}}$  with rational-valued means and mixture coefficients.

We may fix a probability mass function  $w$  on  $\mathcal{P}$ , which we shall think of as the *prior distribution* on  $\mathcal{P}$ . An *estimator* at sample size  $n$  is a function  $\check{p} : \mathcal{Z}^n \rightarrow \mathcal{P}$  that maps each possible sequence of observations  $z^n = (x_1, y_1), \dots, (x_n, y_n) \in \mathcal{Z}^n$  to some  $p \in \mathcal{P}$ . Following e.g. Barron and Cover (1991), the notation  $\check{p}(Z^n)$  denotes the density of the observed data  $Z^n$  according to the  $\check{p}$  that was selected (estimated) based on the same data  $Z^n$ .

**Conditions** We only consider combinations of  $(\mathcal{P}, P^*)$  and prior  $w$  for which **(A)** for each  $(x, y) \in \mathcal{Z}$ , there is a  $p \in \mathcal{P}$  such that  $w(p)p(y | x) > 0$ . We also assume **(B)** that for some finite integer  $L_{\max} > 0$ , for all  $n$ , all  $z^n \in \mathcal{Z}^n$ ,  $-\log w(\check{p}_2) \leq nL_{\max}$ . Hence the codelength of the 2-MDL estimator is of order no larger than  $n$ . This assumption is innocuous, since it can always be satisfied by adding one or a few distributions to  $\mathcal{P}$  (proof sketch in appendix). Finally, let

$$V = V(\mathcal{P}, P^*) := \text{ess sup}_Z \sup_{p, p' \in \mathcal{P}} \frac{p(Z_i)}{p'(Z_i)}. \quad (5)$$

We assume **(C)**  $1 < V$  and **(D)**  $V < \infty$ . We may think of  $V(\mathcal{P}, P^*)$  as the maximum ratio between the density of  $z_i = y | x$  assigned by different  $p \in \mathcal{P}$ , where the maximum is over

all  $(x, y)$  in the support of  $P^*$ . In case  $P^*$  has full support, the essential supremum can be replaced by the standard supremum. Assumptions (A), (B) and (C) are harmless; (D) is further discussed in Section 6.

**Mixing** The safe estimator is the  $\check{\kappa}_{\text{SAFE}}$ -two part estimator, with  $\check{\kappa}_{\text{SAFE}}$  determined by the data. To find  $\check{\kappa}_{\text{SAFE}}$ , we test, for each fixed  $\kappa$ , whether we can get a better fit of the data/additionally compress the data by a convex combination of  $\check{p}_\kappa$  with a single other distribution  $p' \in \mathcal{P}$ . Of course, since  $\mathcal{P}$  may be infinite and arbitrary, it may be the case that, no matter what data we observe, there is *always* some  $p \in \mathcal{P}$  such that a convex combination of  $\check{p}_\kappa$  and  $p$  gives a substantially better fit to the data. This problem, it turns out, can be addressed by only looking at distributions  $p$  with prior mass not much smaller than  $w(\check{p}_\kappa)$ : specifically, we require  $-\log w(p) \leq \lceil -\log w(\check{p}_\kappa) \rceil$ , where  $\lceil \cdot \rceil$  denotes rounding up.

To formalize this, for any  $p, p' \in \mathcal{P}$  and any  $\lambda \in [0, 1]$ , we define the mixture distribution  $\text{MIX}(p, p', \lambda)$  as  $(1 - \lambda)p + \lambda p'$ , so that for a single outcome  $z$ ,  $\text{MIX}(p, p', \lambda)(z_i) := (1 - \lambda)p(z_i) + \lambda p'(z_i)$  (note the somewhat special notation).  $\text{MIX}$  is extended to  $n$  outcomes by independence:

$$\text{MIX}(p, p', \lambda)(z_i^n) := \prod_{i=1}^n ((1 - \lambda)p(z_i) + \lambda p'(z_i)). \quad (6)$$

Our test is defined in terms of how much better fit can be achieved by the best-fitting convex combination of this form. To this end, for an arbitrary estimator  $\check{p}$ , we let

$$\text{SUPMIX}(\check{p})(Z_1^n) := \sup_{p \in \mathcal{P}: -\log w(p) \leq \lceil -\log w(\check{p}) \rceil, \lambda \in [0, 1]} \text{MIX}(p, \check{p}, \lambda)(Z_1^n). \quad (7)$$

Let  $Z(\eta) := \sup_{x \in \mathcal{X}, p \in \mathcal{P}} \int_{y \in \mathcal{Y}} (p(y | x))^\eta dy$ . For each  $p \in \mathcal{P}$ , and each  $\eta \in \mathbb{R}$  with  $Z(\eta) < \infty$ , we define  $p^{(\eta)}(y|x) := (p(y|x))^\eta / Z(\eta)$ . Note that the  $p^{(\eta)}$  all represent distributions which, in general, may be defective, even if  $p$  is not. We define  $\mathcal{P}^{(\eta)} := \{p^{(\eta)} \mid p \in \mathcal{P}\}$ . Since, in all our equations, every occurrence of a density  $p^{(\eta)}$  will actually be as a *ratio* of two densities, i.e.  $p^{(\eta)}(Z_1^n) / q^{(\eta)}(Z_1^n)$ , we can safely write  $p^\eta$  instead of  $p^{(\eta)}$  everywhere without affecting the results. This is what we do below. However, for interpreting our results it is useful to think of the  $p^{(\eta)}$  as densities. For predictor models,  $\eta$  corresponds to  $\beta$  in (2); but see Section 6.

**Convexity Lack** We define the *convexity lack* of an arbitrary estimator  $\check{p}$  on data  $Z_1^n$  as

$$\text{CONV-LACK}(\eta, \check{p}) = -\frac{c_\eta}{\eta} \log \frac{\check{p}^\eta(Z_1^n) w^2(\check{p})}{\text{SUPMIX}(\check{p}^\eta)(Z_1^n)}, \quad (8)$$

where  $c_\eta = 1 + C_\eta C'_\eta$ , and  $C_\eta = 2 + 2\eta \log V$  and  $C'_\eta = 2V^{2\eta}$ . The rationale behind these values will become clear in Theorem 1 and Lemma 9 below.  $\text{CONV-LACK}$  is a measure of how many more bits are needed to encode the data using a two-part code (with  $\kappa = 2$ ) based on  $\check{p}^\eta$  (the numerator in the logarithm) as compared to the number of bits needed by the two-component mixture that provides the best fit (smallest codelength) with hindsight (the denominator). The larger this number, the more could have been gained by modelling the data with the convex hull of  $\mathcal{P}$  rather than just  $\mathcal{P}$ . Our main insight (see Theorem 3 below) is that, if  $\eta \leq \eta_{\text{CRIT}}$  ( $\eta_{\text{CRIT}}$ , introduced in Section 1, is formally defined below), then  $\text{CONV-LACK}$  is guaranteed to be small. This suggests to test various values of  $\kappa$ ,

and define the safe estimator as the  $\kappa$ -two part estimator for the smallest value of  $\kappa$  for which  $\text{CONV-LACK}(1/\kappa, \check{p}_\kappa)$  is below some fixed threshold. Here we opt for the essentially equivalent, but mathematically more convenient option to simply *add*  $\text{CONV-LACK}$  as an additional penalty to the codelength of the two-part estimator:

**Safe Estimation** Let  $\kappa_{\max} = \lceil \sqrt{n}/(2 \log V) \rceil$  (this value is motivated below Lemma 5). The *safe estimator*  $\check{p}_{\text{SAFE}}$  is defined as the  $2\kappa$ -two-part estimator for the  $\kappa \in \{1, 2, \dots, \kappa_{\max}\}$  that minimizes

$$-\log \check{p}_{2\kappa}(Z_1^n) - 2\kappa \log w(\check{p}_{2\kappa}) + \text{CONV-LACK}(\kappa, \check{p}_{2\kappa}). \quad (9)$$

This is just the formula for  $\check{p}_{2\kappa}$ , but with the term  $\text{CONV-LACK}(\kappa, \check{p}_{2\kappa})$  added. Here we establish that  $\check{p}_{\text{SAFE}}$  has good theoretical properties; whether it is useful in practice is discussed in Section 6.

### 3. Generalization Bounds for the Safe Estimator

**Preliminaries** We define the (generalized) *information closure* of  $\mathcal{P}$  (Barron and Cover, 1991) as

$$\langle \mathcal{P} \rangle := \{p' \mid \text{for some } P^*, \inf_{p \in \mathcal{P}} D^*(p^* \| p) = D^*(p^* \| p')\}, \quad (10)$$

where  $p'$  ranges over *all* conditional densities for  $Y$  given  $X$  (i.e.  $p'$  is not necessarily in  $\mathcal{P}$ ), and  $P^*$  ranges over *all* distributions on  $\mathcal{Z}$  that have some conditional density; we denote the density corresponding to  $P^*$  by  $p^*(y \mid x)$ . Henceforth we assume that data  $Z^n$  are sampled from a  $P^*$  that admits such a  $p^*$ . We also assume that  $P^*$  and  $\mathcal{P}$  are such that there exists a *best-approximating*  $q$ , defined as a  $q$  such that:

$$D^*(p^* \| q) = \inf_{p \in \mathcal{P}} D^*(p^* \| p) \text{ and } V(\mathcal{P} \cup \{q\}, P^*) = V(\mathcal{P}, P^*). \quad (11)$$

From now on, for given  $(\mathcal{P}, P^*)$ , we fix a particular best-approximating density once and for all and keep denoting it by  $q$ . We must have  $q \in \langle \mathcal{P} \rangle$ , and  $D^*(q \| p) \geq 0$  for all  $p \in \mathcal{P}$ . Our assumption that  $p^*$  and  $q$  exist simplifies the formulation of our theorems. Still, all our results continue to hold, with appropriately generalized definitions, if no such  $q$  or  $p^*$  exist\*. In the well-specified case, with  $q = p^* \in \langle \mathcal{P} \rangle$ , we trivially have that, for  $\eta = 1$ ,

$$\text{For all } p \in \mathcal{P}: E^* \left[ \left( \frac{p(Z_i)}{q(Z_i)} \right)^\eta \right] \leq 1, \text{ or equivalently, } d_\eta^*(q \| p) := -\frac{1}{\eta} \log E^* \left[ \left( \frac{p(Z_i)}{q(Z_i)} \right)^\eta \right] \geq 0, \quad (12)$$

as is seen by writing out the expectation in full and substituting  $q$  by  $p^*$ . Here  $d_\eta^*$  is the generalized Rényi divergence (Li, 1999); by a result of Li (1999), repeated as Proposition 15 in Section 5, if (12) holds then  $d_{\eta/2}^*$  may be viewed as a proxy for the generalized KL divergence, since then, for all  $p \in \mathcal{P}$ ,  $D^*(q \| p) \leq C_\eta d_{\eta/2}^*(q \| p) \leq C_\eta D^*(q \| p)$ , where  $C_\eta = 2 + 2\eta \log V$  is a constant. This is a key idea for our proofs. Classical theorems on two-part MDL inference for the well-specified case (Barron and Cover, 1991; Zhang, 2006; Grünwald, 2007) invariably make use of (12) at some point in the proofs; so do classical results on Bayesian consistency (Doob, 1949), in which (12) is used to establish that  $\{p(Z_i^n)/q(Z_i^n)\}_{n=1,2,\dots}$  forms a martingale. It can be shown (Li, 1999; Kleijn and van der Vaart, 2006) that (12) still holds for  $\eta = 1$  if  $\langle \mathcal{P} \rangle$  is convex, or, more generally, if (4) holds. This is the fundamental

reason why the MDL and Bayesian convergence bounds still hold in that setting. In fact, (4) with  $\mathcal{P}^{(\eta)}$  in the role of  $\mathcal{P}$  is equivalent to (12), as follows from Lemma 9 in Section 5 (proof sketch in Appendix). If (4) does not hold for  $\eta = 1$ , then (12) does not hold for  $\eta = 1$  and MDL and Bayes may not converge (Example 7 below). Luckily, for many types of  $\mathcal{P}$ , one can still show that (12) holds for some  $\eta < 1$ . Thus it makes sense to define the *critical exponent*  $\eta_{\text{CRIT}}$  as the largest value of  $\eta$  such that, for all  $p \in \mathcal{P}$ , (12) holds. It is useful to generalize the idea slightly. We define, for  $u \geq 0$ , the *u-critical exponent*  $\eta_{\text{CRIT}}(u)$  as

$$\eta_{\text{CRIT}}(u) := \sup \left\{ \eta \leq 1 : \text{for all } p \in \mathcal{P}, \quad \log E^* \left[ \left( \frac{p(Z_i)}{q(Z_i)} \right)^\eta \right] \leq \frac{u}{n} \right\}. \quad (13)$$

$\eta_{\text{CRIT}}(0)$  is just the critical value as defined before. In Section 1 we cheated a little, writing  $\eta_{\text{CRIT}}$  for  $\eta_{\text{CRIT}}(u)$  for the  $u$  which gives the best bounds; see below. Whenever we write “WHP” (‘with high probability’), we really mean “for all  $K \geq 0$ , with  $P^*$ -probability at least  $1 - e^{-K}$ ,  $Z^n$  satisfies...”.

**Theorem 1 (Oracle Bound)** *Let  $(\mathcal{P}, P^*)$  and  $w$  satisfy conditions (A)-(D) of Section 2 with  $V$  as in (5),  $q$  as in (11),  $\eta_{\text{CRIT}}$  as in (13) and  $\kappa_{\text{max}}$  as above (9). Let  $\check{p}$  be an arbitrary estimator. Let  $Z^n \sim P^*$ . WHP, uniformly for all  $\eta \in \{1, 1/2, 1/3, 1/4, \dots, 1/\kappa_{\text{max}}\}$ , all  $u \in \{0, 1, \dots, nL_{\text{max}}\}$ , we have*

$$D^*(q \parallel \check{p}) \leq \frac{C_\eta}{n} \left( \text{RED}(2/\eta, \check{p}) + \text{CONV-LACK}(\eta, \check{p}) + \frac{u}{\min\{\eta, \eta_{\text{CRIT}}(u)\}} + R \right), \quad (14)$$

where  $C_\eta = 2 + 2\eta \log V$ . The term CONV-LACK is given by (8). The term RED is given by

$$\text{RED}(2/\eta, \check{p}) = -\log \frac{w^{2/\eta}(\check{p})\check{p}(Z_i^n)}{q(Z_i^n)} = -\frac{2}{\eta} \log w(\check{p}) - \log \check{p}(Z_i^n) + \log q(Z_i^n). \quad (15)$$

The remainder term\* is  $R = O\left(\frac{K + \log n + \log(2 + \lceil -\log w(\check{p}) \rceil)}{\min\{\eta, \eta_{\text{CRIT}}(u)\}}\right)$ .

RED stands for ‘redundancy’. It can be interpreted as the extra number of nats needed to code the data using a two-part code based on  $\check{p}$  (with  $\kappa = 2/\eta$ ), as compared to the code based on the best-approximating  $q$ , and under mild conditions RED/ $n$  will tend to 0, WHP (Example 4). For predictor models as in (2),  $\check{p} = p_{\check{f}}$  and  $q = p_g$ , and then RED/ $n$  can be thought of as the difference in empirical risk between  $\check{f}$  and the optimal  $g$ , ‘penalized’ by  $-(2/\eta n) \log w(\check{p})$ . The RED and CONV-LACK terms depend on the data; the third term in (14) becomes 0 if we set  $u = 0$ ; for the possibility  $u > 0$  see Example 6. The fourth term is a remainder term which does not depend on the data except for the term  $\log(2 - \log w^{1/\eta}(\check{p}))$  which, by assumption (B), if  $\check{p}$  is a 2-part or the safe estimator, is bounded by  $O(\log n)$ .

We can now motivate the definition  $\check{p}_{\text{SAFE}}$ : among all two-part estimators, it is the one minimizing (14), ignoring the remainder term. To see this, note that the third term in (14) does not depend on the data, and the redundancy term can be written as the sum of two terms plus a term  $\log q(Z_i^n)$  that does not depend on our choice of estimator  $\check{p}$ . Thus,  $\check{p}_{\text{SAFE}}$  minimizes an upper bound on the KL divergence between  $\check{p}$  and the best-approximating  $q$ .

Let us compare (14) to the bounds on the standard two-part MDL estimator as derived by Barron and Cover (1991) under the assumption that  $q = p^*$ . One of their main results



(Theorem 4, as later strengthened by Zhang (2006) and (Grünwald, 2007, Section 15.3)) implies that, for all  $\kappa \geq 1$ ,

$$E_{Z^n \sim P^*} [ D^*(p^* \| \check{p}_{2\kappa}) ] \leq \frac{C}{n} E_{Z^n \sim P^*} [ \text{RED}(2\kappa, \check{p}_{2\kappa}) ], \quad (16)$$

where  $C \leq 2 + 2 \log V$  as above. This provides a frequentist justification for the 2-part MDL estimator  $\check{p}_{2\kappa}$  with  $\kappa \geq 1$ , since  $\check{p}_{2\kappa}$  is in fact defined as the  $p \in \mathcal{P}$  that minimizes  $\text{RED}(2\kappa, p)$ . Apart from the ‘in-expectation’ rather than ‘in-probability’ formulation, the “only” relevant difference to our result is the CONV-LACK term, which appears in (14) because we do not require a correct model. In Section 4 we show that for many combinations of  $P^*$  and  $\mathcal{P}$ , the CONV-LACK term will be small, WHP, for  $\check{p} = \check{p}_{\text{SAFE}}$ . In such cases Theorem 1 states something similar to Barron and Cover’s result ( $D^*(q \| \check{p}) \rightarrow 0$ ), but without the often unrealistic requirement that  $p^* \in \langle \mathcal{P} \rangle$ .

**Theorem 2 (Empirical Bound)** *Assume the notations and conditions of Theorem 1. WHP, uniformly for all  $\eta \in \{1, 1/2, 1/3, \dots, 1/\kappa_{\max}\}$ , we have,*

$$E^*[-\log \check{p}(Z_i)] \leq \frac{1}{n} (-\log \check{p}(Z_i^n) - \log w(\check{p})^{2/\eta} + 2\text{CONV-LACK}(\eta, \check{p}) + R), \quad (17)$$

with remainder term  $R = O\left(\sqrt{n(K + \log n + \log(2 - \log w(\check{p})))}\right)$ .

The proof of this result is similar to the proof of Theorem 1 and will be provided in the full paper. Note that the weights of the main terms on the right side in Theorem 2 and 1 are different. In Theorem 2, the left term’s weight is reduced from  $C_\eta$  to 1, and the weight of the right term (CONV-LACK) is reduced from  $C_\eta$  (which is always  $\geq 2$ ) to 2. Unlike the ‘oracle’ bound Theorem 1, the ‘empirical’ bound Theorem 2 gives useful information without knowledge of  $E^*[-\log q(Z_i)]$ . If  $\mathcal{P} = \mathcal{P}_{\mathcal{F}}$  for a classification model  $\mathcal{F}$ , then the first term on the right-hand side is the empirical risk  $\beta^{-1} n^{-1} \sum_{i=1}^n \text{LOSS}(y_i, \check{f}(x_i))$  and the bound becomes similar to the PAC-Bayesian and Occam’s Razor (OR) bounds (McAllester, 2003; Blumer et al., 1987). Yet, by Condition (B), for  $\check{p} = \check{p}_{\text{SAFE}}$ , the remaining error term  $R/n$  is of order  $O(\log n / \sqrt{n})$  rather than  $O\left(\sqrt{-\log w(\check{p})/n}\right)$  as it would be for PAC-Bayes and OR-bounds. In this sense, for data  $Z^n$  such that for some  $\gamma > 0$  (say,  $\gamma = 0.1$ ), for all  $f \in \mathcal{F}$ , the empirical loss of  $f$  on  $Z^n$  is larger than  $\gamma$ , the bound of Theorem 2 will be stronger than the best PAC-Bayesian or OR bounds. In other cases Theorem 2 gives weaker bounds than PAC-Bayes, since unlike PAC-Bayes it does not improve if  $\check{f}$  has empirical error  $\approx 0$ ; also it is not suitable for randomized classifiers. Theorem 2 is thus a first step, to be improved in future work.

#### 4. What Actually Happens

**Theorem 3** *Assume the notations and conditions of Theorem 1. Let  $c_\eta$  be as in (8). Fix  $u \geq 0$  and let  $\eta \leq \eta_{\text{CRIT}}(u)$ . WHP, we have  $\text{CONV-LACK}(\eta, \check{p}) \leq c_\eta \left(3\text{RED}(2/\eta, \check{p}) + \frac{u}{\eta}\right) + R$ , with remainder term  $R = O((\log n + K + u)/\eta)$*

Applying Theorem 1 to the safe estimator  $\check{p}_{\text{SAFE}}$  with  $\eta \leq \eta_{\text{CRIT}}(u)$ , and using Theorem 3 to rewrite CONV-LACK, and using the fact that, if two inequalities hold with high probability,

the combined inequality also holds with high probability (see Proposition 12 in Section 5), we see that for all  $\eta \in \{1, 1/2, 1/3, \dots, \kappa_{\max}\}$  with  $\eta \leq \eta_{\text{CRIT}}(u)$ , the safe estimator achieves, WHP,

$$\begin{aligned} D^*(q \|\check{p}_{\text{SAFE}}) &\leq \frac{C}{n} \left( \text{RED}(2/\eta, \check{p}_{\text{SAFE}}) + \frac{u}{\eta} + c_\eta 3 \text{RED}(2/\eta, \check{p}_{\text{SAFE}}) + \frac{c_\eta}{\eta} u + R' \right) \\ &= \frac{C(1+3c_\eta)}{n} \left( \text{RED}(2/\eta, \check{p}_{2/\eta}) + R'' \right) + \frac{C(1+c_\eta)}{n} \cdot \frac{u}{\eta} \\ &\leq \frac{C''}{n} \left( \text{RED}(2/\eta, \check{p}_{2/\eta}) + \frac{u}{\eta} + R'' \right), \end{aligned} \quad (18)$$

with constant  $C'' = C(1 + 3c_\eta)$  and new remainder term  $R'' = O((\log n + K)/\eta)$ . As long as we use (18) with  $u = 0$  (directly below) or  $u = 1$  (in Lemma 5) and  $\eta \leq \eta_{\text{CRIT}}(u)$ , then the terms  $u/\eta$  and  $R''$  are at most of the same order as the first term  $(-2/\eta) \log w(\check{p}_{2/\eta})$  in RED, and hence do not affect the obtained convergence rates of  $\check{p}_{\text{SAFE}}$ .

**Example 4 [Best-Case: Model  $\mathcal{P}$  correct or convex]** Suppose that  $P^*$  is in the information closure of  $\mathcal{P}$ , i.e.  $q = p^*$ . Then  $\eta_{\text{CRIT}}(0) = 1$ , and, using (18) with  $u = 0$  and  $\eta = \eta_{\text{CRIT}}(0)$ , WHP,

$$D^*(p^* \|\check{p}_{\text{SAFE}}) = D^*(q \|\check{p}_{\text{SAFE}}) \leq \frac{C''}{n} \left( \text{RED}(2, \check{p}_2) + R'' \right), \quad (19)$$

where by Barron and Cover's original analysis, we would get (16). Except for the in-probability rather than in-expectation formulation, the only real difference is that the KL divergence is bounded in terms of a larger constant factor. This is the price we pay for not knowing in advance that our model was, in fact, correct, while using a procedure that still leads to good results if it is incorrect.

Barron and Cover (1991) show that, for a wide variety of probabilistic models  $\mathcal{M}$ , there exist countable discretizations  $\mathcal{P} \subset \mathcal{M}$  and corresponding priors  $w$  on  $\mathcal{P}$  such that  $\frac{1}{n} E_{Z^n \sim P^*} [\text{RED}(2, \check{p}_2)]$  is equal to the minimax convergence rate in KL risk if  $\mathcal{M}$  is “non-parametric”, or equal to the minimax rate up to a  $\log n$ -factor if  $\mathcal{M}$  is “parametric”. Using Markov's inequality\*, WHP  $\text{RED}(2, \check{p}_2)$  on the data (as in (19)) is not larger than a constant factor times its expectation  $E_{Z^n \sim P^*} [\text{RED}(2, \check{p}_2)]$ . This implies that the standard MDL estimator also achieves the minimax rate in probability (up to a logarithmic factor in the parametric case). Hence, by (19), so do we. A similar story\* can be told if  $\langle \mathcal{P} \rangle$  is convex or if the weaker condition (4) holds; again, up to constant factors, the safe estimator performs as well as the two-part estimator, which converges at near-optimal rates.

**When the Model is Wrong** Define  $D_{\text{SQ}}^*(q, p) = E^*[(\log p(Z_i)/q(Z_i))^2]$ . Such a variation of generalized KL divergence was earlier considered by, e.g., Kleijn and van der Vaart (2006). The lemma below shows that if the model is wrong, then the value of  $\eta_{\text{CRIT}}(u)$  depends on the relation between  $D_{\text{SQ}}^*$  and  $D^*$ . The lemma is not really new, being a direct translation of existing results of e.g. Tsybakov (2004) from ‘ $\mathcal{F}$ -space with loss function LOSS’ to ‘ $\mathcal{P}$ -space with log-loss’.

**Lemma 5** *Assume the notations and conditions of Theorem 1. Suppose further (E) that for some  $A > 0$  and some  $0 \leq \nu \leq 1$ , for all  $p \in \mathcal{P}$ ,  $D_{\text{SQ}}^*(q, p) \leq A(D^*(q \| p))^\nu$ . Then, for all  $u > 0$ , we have*

$$\eta_{\text{CRIT}}(u) \geq \min \left\{ \frac{1}{2 \log V}, B \left( \frac{u}{n} \right)^{\frac{1-\nu}{2-\nu}} \right\} \quad \text{where } B = (2/eA)^{\frac{1}{2-\nu}}. \quad (20)$$

When  $\mathcal{P}$  represents a classification model containing the Bayes classifier, condition (E) above specializes to the celebrated condition of (Mammen and Tsybakov, 1999; Tsybakov, 2004); the  $\kappa$  in (Tsybakov, 2004) is equal to  $\nu^{-1}$  in our notation. In particular, we automatically have  $D_{\text{SQ}}^*(q, p) \leq (\log V)^2$  so (E) always holds for  $\nu = 0$  and  $A = (\log V)^2$ . Using (20) with these values, and using  $1/2 < \sqrt{2/e}$ , it follows that for all  $u \geq 1$ ,

$$\eta_{\text{CRIT}}(u) \geq \frac{1}{2 \log V} \sqrt{\frac{u}{n}} \geq 1/\kappa_{\text{max}}, \quad (21)$$

which explains why we could restrict  $\eta$  to  $\eta > \kappa_{\text{max}}$ ; see Example 6. If (E) holds for some  $\nu > 0$  though, then  $\eta_{\text{CRIT}}(u)$  is of larger order than  $\sqrt{u/n}$  and things get better; see below Example 6.

**Example 6 [Worst-Case]** Let  $u \geq 1$  and let  $\eta = \eta_{\text{CRIT}}(u)$ . Using (18), and (21), we see that WHP,

$$\begin{aligned} D^*(q \| \check{p}_{\text{SAFE}}) &\leq \frac{C''}{n} \left( \text{RED}(2/\eta, \check{p}_{2/\eta}) + \frac{u}{\eta} + R'' \right) \\ &= \frac{C''}{n} \left( -2c\sqrt{\frac{n}{u}} \log w(\check{p}_{2/\eta}) - \log \frac{\check{p}_{2/\eta}(Z_i^n)}{q(Z_i^n)} + c\sqrt{u \cdot n} + R'' \right) \end{aligned}$$

for some constant  $c = 2 \log V$ . Differentiating with respect to  $u$  shows that a minimum is achieved\* for  $u \approx -2 \log w(\check{p}_{2/\eta})$ . The resulting expression becomes

$$C''' \cdot 2c\sqrt{\frac{-\log w(\check{p}_{2/\eta})}{n}} + C''' \left( -\frac{1}{n} \log \frac{\check{p}_{2/\eta}(Z_i^n)}{q(Z_i^n)} + R'' \right). \quad (22)$$

For classification models, this bound is familiar from the computational learning literature. Now suppose that (E) holds for some  $\nu > 0$ . Then we can achieve better bounds: by the reasoning below Theorem 3,  $\check{p}_{\text{SAFE}}$  converges at the same rate as the  $\kappa$ -MDL estimator with  $\kappa = \eta_{\text{CRIT}}(1)^{-1} = O(n^{\frac{1-\nu}{2-\nu}})$ . From (18) we see that in the special case that  $q \in \mathcal{P}$ ,  $w(q) > 0$ , (18) gives a rate in probability of  $(\log n)/n^{1/(2-\nu)}$ , which, for classification models, is equal to the minimax optimal rate in expectation (Tsybakov, 2004) up to a log factor. The next example illustrates this for  $\nu = 1$ ; in the full paper\* we will also provide examples involving regression and classification with  $0 < \nu < 1$ .

**Example 7 [Bayesian inconsistency and Tsybakov's Condition]**

Grünwald and Langford (2007) showed that standard MDL and Bayesian inference can be inconsistent in various ways if  $P^* \notin \langle \mathcal{P} \rangle$ , for countable models  $\mathcal{P} = \{p_0, p_1, \dots\}$  that are really classification models, i.e.  $\mathcal{P} = \mathcal{P}_{\mathcal{F}}$  with  $\mathcal{F} = \{f_0, f_1, \dots\}$  with  $p_j = p_{f_j}$  as given by (2), where  $\mathcal{Y} = \{0, 1\}$  and LOSS is the 0/1-LOSS. In these examples,  $p_0$  has positive prior  $w(p_0) > 0$  independent of the sample size, and for some  $\delta > 0$ , for all  $j > 0$ , it holds  $D^*(p_0 \| p_j) > \delta$ , i.e.  $\text{RISK}(p_j) > \text{RISK}(p_0) + \beta^{-1}\delta$ . In the examples Tsybakov's condition (E) holds with  $\nu = 1$  but only for very large  $A$ . Since thus  $q = p_0$  and, by Lemma 5,  $\eta_{\text{CRIT}}(1) > 1/C$  for some very large constant independent of  $n$ , it follows from (18) that the safe estimator converges WHP at rate  $O((-\log w(p_0) + R'')/n) = O(\log n/n)$ , much faster than the worst-case  $O(1/\sqrt{n})$ . However, explicit calculation of  $\eta_{\text{CRIT}}(1)$  shows that it is indeed very small, and since standard MDL and Bayesian MAP use an  $\eta$  equal to 1 or 2, it comes as no surprise that in this scenario they do not converge at all, i.e. with probability 1, for all large  $n$ , they select a distribution/classifier  $p \neq p_0$ , as was shown formally by Grünwald and Langford (2007).

## 5. The Proofs

**Preliminary Results** Our main tool is Proposition 8 below, a bound for ratios of probability densities, similar to earlier inequalities by Barron and Cover (1991); Li (1999); Zhang (2006). Below  $\text{TR}$  is a function mapping distributions to other distributions (we use notation as in (6)).

**Proposition 8** *Let  $Z^n$  be i.i.d.  $\sim P^*$ . Let  $\mathcal{P}$  be a countable set of (possibly defective) conditional densities for  $\mathcal{Z}^n$  and let  $\check{p}$  be an arbitrary estimator. Let  $\mathcal{Q}$  be another set of (possibly defective) conditional densities for  $\mathcal{Z}^n$ . Let  $\text{TR} : \mathcal{P} \rightarrow \mathcal{Q}$  be a function mapping distributions in  $\mathcal{P}$  to distributions in  $\mathcal{Q}$ . Let  $w$  be a (potentially defective) probability mass function on  $\mathcal{P}$ . Let  $\eta > 0$ . Then WHP,*

$$d_\eta^*(\text{TR}(\check{p})\|\check{p}) \leq \frac{1}{n} \left( -\log \frac{w(\check{p})\check{p}(Z_1^n)}{\text{TR}(\check{p})(Z_1^n)} + \frac{K}{\eta} + \frac{1}{\eta} \log \sum_{p \in \check{\mathcal{P}}} w(p)^\eta \right). \quad (23)$$

**Proof** We bound the probability that (23) does *not* hold:

$$\begin{aligned} & P^* \left( \eta d_\eta^*(\text{TR}(\check{p})\|\check{p}) > \frac{1}{n} \left( -\eta \log \frac{w(\check{p})\check{p}(Z_1^n)}{\text{TR}(\check{p})(Z_1^n)} + K + \log \sum_{p \in \check{\mathcal{P}}} w(p)^\eta \right) \right) = \\ & P^* \left( \log \left( \frac{\check{p}(Z_1^n)}{\text{TR}(\check{p})(Z_1^n)} \right)^\eta > -\log \frac{w^\eta(\check{p})}{\sum w^\eta(p)} + K + n \log E^* \left( \frac{\check{p}(Z_1)}{\text{TR}(\check{p})(Z_1)} \right)^\eta \right) \leq \\ & P^* \left( \text{There exists } p \in \mathcal{P} \text{ with } \left( \frac{p(Z_1^n)}{\text{TR}(p)(Z_1^n)} \right)^\eta > e^K \left( \frac{\sum w(p)^\eta}{w^\eta(p)} \right) E^* \left( \frac{p(Z_1^n)}{\text{TR}(p)(Z_1^n)} \right)^\eta \right) \leq \\ & \sum_{p \in \mathcal{P}} P^* \left( \left( \frac{p(Z_1^n)}{\text{TR}(p)(Z_1^n)} \right)^\eta > e^K \left( \frac{\sum w(p)^\eta}{w^\eta(p)} \right) E^* \left( \frac{p(Z_1^n)}{\text{TR}(p)(Z_1^n)} \right)^\eta \right) \leq e^{-K}, \end{aligned}$$

where the equality is basic rewriting, the first inequality follows from exponentiating both sides, absorbing  $n$  into the expectation  $E^*$  (which can be done since the  $Z_i$  are i.i.d.) and weakening, the second is the union bound, and the third is an instance of Markov's inequality.  $\blacksquare$

In some applications we set, for all  $p \in \mathcal{P}$ ,  $\text{TR}(p)$  equal to the best approximating density  $q$ , and then the first term on the right in (23) is equal to  $\text{RED}(1/\eta, \check{p})$ ; the inequality is then a weakening of Zhang's, who provides an expectation rather than an in-probability form. In other applications (e.g. below Eq. (37)),  $\text{TR}(p)$  actually varies with  $p$ , and in this form, the inequality is new. We will apply this proposition in two different ways. In the first type of application, the goal is to get a (high-probability) upper bound on the left-hand side of (23). In the second type of application, the goal is to upper bound  $-\log \text{TR}(\check{p})(Z_1^n)$ . Thus, we rewrite (23) equivalently as:

$$-\frac{1}{n} \log \text{TR}(\check{p})(Z_1^n) \leq -\frac{1}{n} \log w(\check{p})\check{p}(Z_1^n) + R, \text{ with } R = -d_\eta^*(\text{TR}(\check{p})\|\check{p}) + \frac{K + \log \sum_{p \in \check{\mathcal{P}}} w(p)^\eta}{n\eta}. \quad (24)$$

In such applications, we take a value of  $\eta$  guaranteeing  $d_\eta^*(\text{TR}(\check{p})\|\check{p}) \geq 0$  or (as e.g. in Lemma 14), we allow  $d_\eta^*(\text{TR}(\check{p})\|\check{p})$  to be negative but not too negative, so that the bound remains useful.

Let  $p, p' \in \mathcal{P}$  such that  $D^*(p\|p') \geq 0$ , and let  $\lambda^\circ := \arg \min_{\lambda \in [0,1]} D^*(p^*\|\text{MIX}(p, p', \lambda)) = \arg \max_{\lambda \in [0,1]} D^*(\text{MIX}(p, p', \lambda)\|p)$  (if more than one  $\lambda$  achieves the extremum, we take the smallest). Our second key result states that if  $D^*(\text{MIX}(p, p', \lambda^\circ)\|p)$  is small, then  $E^*[\text{MIX}(p, p', \lambda^\circ)(Z_1)/p(Z_1)]$  cannot be much larger than 1; equivalently  $d_1^*(p\|\text{MIX}(p, p', \lambda^\circ))$  cannot be much smaller than 0:

**Lemma 9** Let  $(\mathcal{P}, P^*)$  be as on page 401,  $V = V(\mathcal{P}, P^*)$  be as in (5), and let  $p, p' \in \mathcal{P}$  be such that  $D^*(p||p') \geq 0$ , and  $\lambda^\circ$  be as above. (a) If  $\lambda^\circ = 0$  ( $p$  is closer to  $p^*$  than any mixture of  $p$  and  $p'$ ) then for all  $\lambda \in [0, 1]$ ,  $d_1^*(p||\text{MIX}(p, p', \lambda)) \geq 0$ ; otherwise, (b),  $-d_1^*(p||\text{MIX}(p, p', \lambda^\circ)) \leq 2V^2 D^*(\text{MIX}(p, p', \lambda^\circ)||p)$ .

**Proof** Let  $g(\lambda) = D^*(\text{MIX}(p, p', \lambda)||p)$ . Then  $g(0) = 0, g(1) \leq 0$ . We first need the following (proof straightforward by differentiation, see appendix):

**Proposition 10** 1.  $g'(0) = E^*\left(\frac{p'(Z_i)}{p(Z_i)}\right) - 1$ ; if  $\lambda^\circ = 0$  then (1a)  $g'(0) \leq 0$ ; if  $\lambda^\circ > 0$  then (1b)  $g'(0) > 0, g'(\lambda^\circ) = 0$  and  $g'(1) \leq 0$ ; 2. if  $p(Z_i) = p'(Z_i)$   $P^*$ -almost surely, then (2a)  $g'(\lambda) = g''(\lambda) = 0$  on  $\lambda \in [0, 1]$ . Otherwise (2b)  $g''(\lambda) < 0$  on  $[0, 1]$  and  $\max_{\lambda \in [0, 1]} |g''(\lambda)| \leq \min_{\lambda \in [0, 1]} V^2 |g''(\lambda)|$ .

Abbreviate  $d_1^*(p||\text{MIX}(p, p', \lambda))$  to  $d^*(\lambda)$  and  $g(\lambda^\circ) = D^*(\text{MIX}(p, p', \lambda^\circ)||p)$  to  $D^*$ , and note that

$$-d^*(\lambda) = \log E^*\left(\frac{(1-\lambda)p(Z_i) + \lambda p'(Z_i)}{p(Z_i)}\right) \leq E^*\left(\frac{(1-\lambda)p(Z_i) + \lambda p'(Z_i)}{p(Z_i)}\right) - 1 = \lambda g'(0). \quad (25)$$

In case (1a) and (2a) the result is now immediate, so assume (1b) and (2b). Then by a first-order Taylor approximation of  $g'$ , for some  $0 \leq \lambda_1 \leq \lambda^\circ$ ,  $g'(0) = g'(\lambda^\circ) - \lambda^\circ g''(\lambda_1) = \lambda^\circ |g''(\lambda_1)|$ , so that (25) gives  $-d^*(\lambda^\circ) \leq (\lambda^\circ)^2 |g''(\lambda_1)|$ . Also, by a 2nd-order Taylor expansion of  $g$  around  $\lambda^\circ$  we find, for some  $0 \leq \lambda_2 \leq \lambda^\circ$ , that  $0 = g(0) = (1/2)(\lambda^\circ)^2 g''(\lambda_2) + g(\lambda^\circ)$ , so  $D^* = (1/2)(\lambda^\circ)^2 |g''(\lambda_2)|$ . Combining with our expression for  $d^*(\lambda^\circ)$ , we get  $\frac{-d^*(\lambda^\circ)}{D^*} \leq 2 \frac{|g''(\lambda_1)|}{|g''(\lambda_2)|}$ . The result now follows by part (2b) of Proposition 10. ■

The next proposition is about varying exponents rather than mixture coefficients:

**Proposition 11** Let  $\mathcal{P} = \{p, p'\}$  be such that  $V(\mathcal{P}, P^*) < \infty$  and  $D^*(p||p') > 0$ . Then (a): letting  $g(\eta) = \log E^*(p'(Z_i)/p(Z_i))^\eta = -\eta d_\eta^*(p||p')$ , we have  $g(0) = 0$ ,  $g(\eta)$  is decreasing at  $\eta = 0$  and  $\exp(g(\eta))$  is strictly convex, so that there exists at most one  $\eta' > 0$  with  $g(\eta') = 0$ , and  $g(\eta)$  is increasing for  $\eta \geq \eta'$ . (b) Define  $\lambda^\circ$  as in Lemma 9. If  $\lambda^\circ = 0$  ( $p$  is closer to  $p^*$  than any mixture of  $p$  and  $p'$ ) then  $\forall \eta \in (0, 1)$ ,  $d_\eta^*(p||p') > 0$ .

**Proof** (a) is just differentiation of  $E^*(p'(Z_i)/p(Z_i))^\eta$  (see proof of Lemma 5); details omitted. (b) is immediate from (a) because by Lemma 9, part (a),  $d_1^*(p||p') = d_1^*(p||\text{MIX}(p, p', 1)) = -g(1) \geq 0$ . ■

The following proposition provides the glue that ties all our inequalities together:

**Proposition 12 (log-Bonferroni)** Let  $\mathcal{J}$  be a finite or countably infinite set. Let  $\{Y_j\}_{j \in \mathcal{J}}$  be a collection of random variables, let  $\{a_j\}_{j \in \mathcal{J}}$  be a collection of constants in  $\mathbb{R}$  and let  $\{f_j\}_{j \in \mathcal{J}}$  be a collection of increasing functions  $\mathbb{R} \rightarrow \mathbb{R}$ . Suppose that for all  $j \in \mathcal{J}$ , WHP,  $Y_j \leq a_j + f_j(K)$ . Then, for any collection of positive numbers  $\{w_j\}_{j \in \mathcal{J}}$  such that  $\sum_{j \in \mathcal{J}} w_j = 1$ , we have, WHP,

$$\text{For all } j \in \mathcal{J}, Y_j \leq a_j + f_j(K - \log w_j).$$

This result is a straightforward consequence of the union bound that appears in one form or other in many COLT papers; for convenience there is a proof in the appendix.

**Notation common to the Proofs** In all proofs below we make use of the following concepts: let  $w$  be a prior for a countable set of densities  $\mathcal{P}$ . Let  $p \in \mathcal{P}$ . Relative to  $w$  and  $P^*$ , we define\* the *optimal density at  $p$ 's description length* as

$$\text{OPT}(p) := \arg \min_{p' \in \mathcal{P}: -\log w(p') \leq \lceil -\log w(p) \rceil} D^*(p^* \| p) \quad (26)$$

For an estimator  $\check{p}$ ,  $\text{OPT}(\check{p})$  is itself a random variable, representing the best distribution (closest in KL divergence to  $p^*$ ) with prior no smaller (up to rounding) (or “complexity” no larger, up to rounding) than the  $\check{p}$  selected for the given data  $Z_1^n$ . We further define

$$\text{OPT}(\mathcal{P}) = \{p \in \mathcal{P} : p = \text{OPT}(p') \text{ for some } p' \in \mathcal{P}\}. \quad (27)$$

Now define for  $p \in \mathcal{P}$ ,  $\text{OPTMIX}(p) := \text{MIX}(\text{OPT}(p), p, \lambda^\circ)$ , where  $\lambda^\circ \in [0, 1]$  minimizes\*

$$E^*[-\log \text{MIX}(\text{OPT}(p), p, \lambda)(Z_1)] = E^*[-\log((1 - \lambda)\text{OPT}(p)(Z_1) + \lambda p(Z_1))]. \quad (28)$$

Note that  $D^*(\text{OPTMIX}(\check{p}^\eta) \| \check{p}^\eta) = \max_{\lambda \in [0, 1]} D^*(\text{MIX}(\text{OPT}(\check{p}^\eta), \check{p}^\eta, \lambda) \| \check{p}^\eta) \geq D^*(\text{OPT}(\check{p}^\eta) \| \check{p}^\eta)$ . (29)

### 5.1. Proofs of Main Results

**Proof of Theorem 1** We first consider an arbitrary fixed  $\eta$  and a fixed  $u$ . We have:

$$D^*(q \| \check{p}) = D^*(q \| \text{OPT}(\check{p})) + \frac{1}{\eta} D^*(\text{OPT}(\check{p}^\eta) \| \check{p}^\eta) \leq D^*(q \| \text{OPT}(\check{p})) + \frac{1}{\eta} D^*(\text{OPTMIX}(\check{p}^\eta) \| \check{p}^\eta), \quad (30)$$

where the equality is straightforward from the definition of  $D^*$  and the inequality follows from (29). By Lemma 14, we can bound the term  $D^*(q \| \text{OPT}(\check{p}))$  from above and rewrite (30) to get, WHP,

$$D^*(q \| \check{p}) \leq \frac{C_\eta}{n} \left( -\frac{1}{\eta} \log \frac{w(\check{p}) \text{OPTMIX}(\check{p}^\eta)(Z_1^n)}{q^\eta(Z_1^n)} + \frac{u}{\eta'} + R_1 \right) + T \quad (31)$$

with  $C_\eta = 2 + 2\eta \log V$ ,  $\eta' = \min\{\eta, \eta_{\text{CRIT}}(u)\}$ ,  $R_1$  as in Lemma 14 and

$$T = -\frac{C_\eta}{\eta} d_1^*(\text{OPT}(\check{p}^\eta) \| \text{OPTMIX}(\check{p}^\eta)) + \frac{1}{\eta} D^*(\text{OPTMIX}(\check{p}^\eta) \| \check{p}^\eta), \quad (32)$$

We proceed to rewrite  $T$ .  $-d_1^*(\text{OPT}(\check{p}^\eta) \| \text{OPTMIX}(\check{p}^\eta))$  may very well be *positive*, but by Lemma 9, applied with  $(\mathcal{P}, p, p') \leftarrow (\mathcal{P}^{(\eta)}, \text{OPT}(\check{p})^\eta, \check{p}^\eta)$  (the notation indicates that e.g.  $\mathcal{P}$  in Lemma 9 is instantiated to  $\mathcal{P}^{(\eta)}$  as used above), we can bound (32) to get:  $T \leq \eta^{-1}(C_\eta C'_\eta + 1) D^*(\text{OPTMIX}(\check{p}^\eta) \| \check{p}^\eta)$ , where we used  $D^*(\text{OPTMIX}(\check{p}^\eta) \| \text{OPT}(\check{p})^\eta) \leq D^*(\text{OPTMIX}(\check{p}^\eta) \| \check{p}^\eta)$ , and  $C'_\eta = 2V^{2\eta}$ . Letting  $w'(p^\eta) := w(p)$ , this can be further bounded using Lemma 13 below, with  $(\check{p}, \mathcal{P}, w) \leftarrow (\check{p}^\eta, \mathcal{P}^{(\eta)}, w')$  (notation as explained above). This gives, WHP:

$$T \leq \frac{C_\eta}{n} \left( -\frac{c_\eta}{\eta} \log \frac{\check{p}^\eta(Z_1^n) w(\check{p})^2}{\text{SUPMIX}(\check{p}^\eta)(Z_1^n)} + R'_1 \right) = \frac{C_\eta}{n} (\text{CONV-LACK}(\eta, \check{p}) + R'_1), \quad (33)$$

where  $R'_1 = \eta^{-1}c_\eta 2K$  and  $c_\eta = C_\eta C'_\eta + 1$  and we used the definition of CONV-LACK. Now apply Proposition 8 as in (24), with  $(\mathcal{P}, \check{p}, \mathcal{Q}, \text{TR}(\cdot), w, \eta) \leftarrow (\mathcal{P}^{(\eta)}, \check{p}^\eta, \mathcal{Q}, \text{OPTMIX}(\cdot), w', 1)$ , where  $\mathcal{Q} = \{\text{OPTMIX}(p^\eta) \mid p \in \mathcal{P}\}$  and  $w'$  as above. Since, by Lemma 9, part (b),  $d_1^*(\text{OPTMIX}(p^\eta) \parallel \text{OPT}(p^\eta)) \geq 0$ , we find that WHP,  $-\log \text{OPTMIX}(\check{p}^\eta)(Z_1^n) \leq -\log w(\check{p})\check{p}^\eta(Z_1^n) + K$ . Substituting this and (33) into (31), we get with Proposition 12 (with  $|\mathcal{J}| = 3, w_j = 1/3$ ) that WHP,

$$D^*(q \parallel \check{p}) \leq \frac{C_\eta}{n} \left( -\frac{1}{\eta} \log \frac{w(\check{p})^2 \check{p}^\eta(Z_1^n)}{q^\eta(Z_1^n)} + \frac{u}{\eta'} + \text{CONV-LACK}(\eta, \check{p}) + R_2 + R'_2 \right), \quad (34)$$

where  $R_2 = (4(K + \log 3) + 4 \log(2 + \lceil -\log w(\check{p}) \rceil)) / \eta'$  and  $R'_2 = c_\eta(2(K + \log 3)) / \eta$ . The result now follows for a fixed value of  $u$  and  $\eta$ . To prove that it holds uniformly for  $u \in \{0, 1, 2, \dots, nL_{\max}\}$ ,  $\eta \in \{1, 1/2, 1/3, \dots, 1/\kappa_{\max}\}$ , we use Proposition 12; details omitted\*.

**Lemma 13** *Let  $(\mathcal{P}, P^*)$  and  $w$  be as on page 401. We have WHP,*

$$D^*(\text{OPTMIX}(\check{p}) \parallel \check{p}) \leq \frac{C_1}{n} \left( -\log \frac{\check{p}(Z_1^n) w^2(\check{p})}{\text{SUPMIX}(\check{p})(Z_1^n)} + 2K \right)$$

where  $C_1 = 2 + 2 \log V$  is a constant and SUPMIX is defined as in (7) above.

**Proof** By Proposition 15, applied with  $\eta = 1$ , we get

$$D^*(\text{OPTMIX}(\check{p}) \parallel \check{p}) \leq C_1 d_{1/2}^*(\text{OPTMIX}(\check{p}) \parallel \check{p}) - (C_1 - 1) d_1^*(\text{OPTMIX}(\check{p}) \parallel \check{p}) \leq C_1 d_{1/2}^*(\text{OPTMIX}(\check{p}) \parallel \check{p}),$$

where  $C_1 = 2 + 2 \log V$  and the final inequality follows because by Proposition 11,  $d_1^*(\text{OPTMIX}(\check{p}) \parallel \check{p}) \geq 0$ . We now let  $\mathcal{Q} = \{\text{MIX}(p_0, p_1, \lambda) : p_0, p_1 \in \mathcal{P}, \lambda \in [0, 1]\}$  be the set of two-component mixtures of elements of  $\mathcal{P}$ , and apply Proposition 8, with  $(\mathcal{P}, \check{p}, \mathcal{Q}, \text{TR}(\cdot), w, \eta) \leftarrow (\mathcal{P}, \check{p}, \mathcal{Q}, \text{OPTMIX}(\cdot), w^2, 1/2)$  (notation as below (32)). We get that, WHP,  $d_{1/2}^*(\text{OPTMIX}(\check{p}) \parallel \check{p}) \leq -\frac{1}{n} \log \frac{w^2(\check{p})\check{p}(Z_1^n)}{\text{OPTMIX}(\check{p})(Z_1^n)} + \frac{2K}{n} \leq -\frac{1}{n} \log \frac{w^2(\check{p})\check{p}(Z_1^n)}{\text{SUPMIX}(\check{p})(Z_1^n)} + \frac{2K}{n}$ . ■

**Lemma 14** *Assume the conditions and notation of Theorem 1. For all  $0 < \eta \leq 1$ , WHP,*

$$D^*(q \parallel \text{OPT}(\check{p})) \leq \frac{C_\eta}{n} \left( -\frac{1}{\eta} \log \frac{w(\check{p}) \text{OPTMIX}(\check{p}^\eta)(Z_1^n)}{q^\eta(Z_1^n)} + \frac{u}{\eta'} + R_1 \right) - \frac{C_\eta}{\eta} d_1^*(\text{OPT}(\check{p}^\eta) \parallel \text{OPTMIX}(\check{p}^\eta)),$$

where  $C_\eta = 2 + 2\eta \log V$ ,  $\eta' = \min\{\eta, \eta_{\text{CRIT}}(u)\}$  and remainder  $R_1 = (3K + 4 \log(2 + \lceil -\log w(\check{p}) \rceil)) / \eta'$ . (Note that  $d_1^*(\text{OPT}(\check{p}^\eta) \parallel \text{OPTMIX}(\check{p}^\eta))$  may be negative).

**Proof** We apply Proposition 15 with  $\eta$  set to  $\eta'$ . This gives  $D^*(q \parallel \text{OPT}(\check{p})) \leq C_{\eta'} d_{\eta'/2}^*(q \parallel \text{OPT}(\check{p})) + (C_{\eta'} - 1)R$ , where  $R = -d_{\eta'}^*(q \parallel \text{OPT}(\check{p}))$ . By Proposition 11, part(a), if  $\eta'R \geq 0$  then, since  $\eta' \leq \eta_{\text{CRIT}}(u)$ , we have  $\eta'R \leq -\eta_{\text{CRIT}}(u) d_{\eta_{\text{CRIT}}(u)}^*(q \parallel \text{OPT}(\check{p})) \leq u/n$ . It follows that  $R \leq u/(\eta'n)$ , so we have

$$D^*(q \parallel \text{OPT}(\check{p})) \leq C_{\eta'} d_{\eta'/2}^*(q \parallel \text{OPT}(\check{p})) + \frac{1}{\eta'} (C_{\eta'} - 1) \frac{u}{n}. \quad (35)$$

Now fix some  $p_0 \in \mathcal{P}$ . Set  $w'(p_0) = 1$ ,  $\text{TR}(p_0) = q$  and apply Proposition 8 with  $(\mathcal{P}, \check{p}, \mathcal{Q}, \text{TR}(\cdot), w, \eta) \leftarrow (\mathcal{P}, p_0, \{q\}, \text{TR}(\cdot), \eta'/2)$ . The  $\check{p}$  in the proposition is a degenerate

estimator that is always equal to the fixed  $p_0$  and does not depend on the data, and  $\text{TR}(\check{p})$  is always equal to  $q$ . We get that WHP,

$$d_{\eta'/2}^*(q\|p_0) \leq \frac{1}{n} \left( -\log \frac{p_0(Z_1^n)}{q(Z_1^n)} + \frac{2K}{\eta'} \right). \quad (36)$$

Note that we can enumerate the elements of  $\text{OPT}(\mathcal{P})$  as given by (27) as  $\{p_1, p_2, \dots\}$  where, for all  $j$ , it holds  $j-1 \leq \lceil -\log w(p_j) \rceil \leq j$ . Using the log-Bonferroni Proposition 12 with  $|\mathcal{J}| = \mathbb{N}$ ,  $w_j = 1/j(j+1)$  and  $f_j(K) = 2K/\eta'(u)$ , we get that, WHP, uniformly for all  $p_j \in \text{OPT}(\mathcal{P})$ ,

$$\begin{aligned} d_{\eta'/2}^*(q\|p_j) &\leq \frac{1}{n} \left( -\log \frac{p_j(Z_1^n)}{q(Z_1^n)} + \frac{2}{\eta'} (K + 2 \log(j+1)) \right) \\ &\leq \frac{1}{n} \left( -\log \frac{p_j(Z_1^n)}{q(Z_1^n)} + \frac{2}{\eta'} (K + 2 \log(2 - \lceil -\log w(p_j) \rceil)) \right). \end{aligned} \quad (37)$$

We now let  $\mathcal{R} = \{\text{OPTMIX}(p^\eta) \mid p \in \mathcal{P}\}$ , and define a prior  $w'$  on  $\mathcal{R}$  with, for  $r \in \mathcal{R}$ ,  $w'(r) := w(p)$  for the  $p$  such that  $r = \text{OPTMIX}(p^\eta)$ . We set  $\text{TR}(p') := \text{OPT}(\check{p}^\eta)$ . We now use Proposition 8 again (in the form (24)) with  $(\mathcal{P}, \check{p}, \mathcal{Q}, \text{TR}(\cdot), w, \eta) \leftarrow (\mathcal{R}, \text{OPTMIX}(\check{p}^\eta), \text{OPT}(\mathcal{P})^\eta, \text{TR}(\cdot), w', 1)$ . (notation as below (32); effectively we use  $\text{OPTMIX}(\check{p}^\eta)$  as an estimator here. ). We get, WHP,

$$-\frac{1}{n} \log \text{OPT}(\check{p}^\eta)(Z_1^n) \leq -\frac{1}{n} \log w(\check{p}) \text{OPTMIX}(\check{p}^\eta)(Z_1^n) - d_1^*(\text{OPT}(\check{p}^\eta) \|\text{OPTMIX}(\check{p}^\eta)) + \frac{K}{n}. \quad (38)$$

Dividing (38) by  $\eta$ , using  $\eta^{-1} \log \text{OPT}(\check{p}^\eta) = \log \text{OPT}(\check{p})$ , and then combining with (35) and (37), with  $p_j$  set to  $\text{OPT}(\check{p})$ , we get, WHP,

$$\begin{aligned} D^*(q\|\text{OPT}(\check{p})) &\leq C_{\eta'} \\ &\left( \frac{1}{n} \left( -\frac{1}{\eta} \log \frac{w(\check{p}) \text{OPTMIX}(\check{p}^\eta)(Z_1^n)}{q^\eta(Z_1^n)} + R \right) - \frac{1}{\eta} d_1^*(\text{OPT}(\check{p}^\eta) \|\text{OPTMIX}(\check{p}^\eta)) + \frac{1}{\eta'} \frac{u}{n} \right), \end{aligned} \quad (39)$$

where  $R = (2/\eta')(K + 2 \log(2 - \lceil -\log w(p_j) \rceil)) + (1/\eta)K$ , which is no greater than  $R_1$  in the statement of the lemma. This proves the result for  $\eta \leq \eta_{\text{CRIT}}(u)$  (for then  $\eta' = \eta$ ). For the case that  $\eta > \eta_{\text{CRIT}}(u)$ , note that then  $C_\eta > C_{\eta'}$ . Because by definition of  $q$  the left-hand side of (39) must be nonnegative, we have WHP that both (39) holds and its right-hand side is nonnegative, so that with the same probability, (39) holds with  $C_{\eta'}$  replaced by  $C_\eta$ . The result follows.  $\blacksquare$

**Proposition 15** *Let  $P^*$  be a distribution on  $\mathcal{Z}$ , let  $p$  and  $q$  be conditional distributions for  $Y$  given  $X$ , and let  $V = V(\{p, q\}, P^*)$  defined as in (5). For all  $\eta \leq 1$  and all  $C_\eta \geq 2 + 2\eta \log V$ , we have*

$$D^*(q\|p) \leq C_\eta \cdot d_{\eta/2}^*(q\|p) - (C_\eta - 1)d_\eta^*(q\|p).$$

*In particular, if  $d_\eta^*(q\|p) \geq 0$ , then  $D^*(q\|p) \leq C_\eta d_{\eta/2}^*(q\|p)$ , i.e. the generalized KL divergence is upper bounded by a constant times the generalized Rényi divergence of order  $1/2\eta$ .*



**Proof** This result is a straightforward extension of a result due to Andrew Barron and Jonathan Li, published in Li's (1999) thesis. See Lemma 5.11, page 67 and Lemma 5.12, page 73 of (Li, 1999).  $d_\eta^*$  corresponds to  $\log c = \log \int fg/f^*$  in Li's notation;  $p^*$  corresponds to  $f$  in Li's notation,  $p(\cdot | x)^\eta$  corresponds to  $g$  in Li, and  $q(\cdot | x)^\eta$  corresponds to  $f^*$ . Our argument is slightly more involved than Li's since we allow conditioning on  $x$ ; this also accounts for the extra  $\eta \log V$  term in the constant  $C_\eta$ . For convenience, we provide a full proof in the appendix.  $\blacksquare$

**Proof of Theorem 3** We fix a  $\hat{p}$  and  $\hat{\lambda}$  be such that  $\text{SUPMIX}(\check{p}^\eta) = \text{MIX}(\hat{p}^\eta, \check{p}^\eta, \hat{\lambda})$ , i.e.  $\hat{p}$  and  $\hat{\lambda}$  achieve\* the supremum in (7) applied with estimator  $\check{p}^\eta$ .

Let  $\hat{\lambda} = \arg \min_{\lambda \in \{0, 1/n, 2/n, \dots, 1\}} -\log \text{MIX}(\hat{p}^\eta, \check{p}^\eta, \lambda)(Z_1^n)$  be a discretized version of  $\hat{\lambda}$ . We have:

$$\begin{aligned} -\log \frac{\check{p}^\eta(Z_1^n)w^2(\check{p})}{\text{SUPMIX}(\check{p}^\eta)(Z_1^n)} &\leq -\log \frac{q^\eta(Z_1^n)}{\text{SUPMIX}(\check{p}^\eta)(Z_1^n)} - \log \frac{\check{p}^\eta(Z_1^n)w^2(\check{p})}{q^\eta(Z_1^n)} \\ &\leq -\log \frac{q^\eta(Z_1^n)}{\text{MIX}(\hat{p}^\eta, \check{p}^\eta, \hat{\lambda})(Z_1^n)} + V^\eta + \eta \text{RED}(2/\eta, \check{p}), \end{aligned} \quad (40)$$

where in the second inequality we used a simple first-order Taylor approximation, showing that  $-\log \text{MIX}(\hat{p}^\eta, \check{p}^\eta, \hat{\lambda})(Z_1^n) \leq -\log \text{MIX}(\hat{p}^\eta, \check{p}^\eta, \hat{\lambda})(Z_1^n) + V^\eta$  (details omitted). We now come to the crucial step: we will prove that WHP, we have

$$-\log \frac{q^\eta(Z_1^n)}{\text{MIX}(\hat{p}^\eta, \check{p}^\eta, \hat{\lambda})(Z_1^n)} \leq -\log w^2(\check{p}) + \log(n+1) + K + u. \quad (41)$$

This result follows by applying Proposition 8 to an extended model  $\mathcal{P}'$  with a prior  $w'$  defined, at sample size  $n$ , as follows:  $\mathcal{P}' = \{\text{MIX}(p_0^\eta, p_1^\eta, \lambda) : p_0^\eta, p_1^\eta \in \mathcal{P}, \lambda \in [0, 1]\}$ ; and, for  $\lambda \in \Lambda := \{0, 1/n, \dots, 1\}$ ,  $w'(\text{MIX}(p_0^\eta, p_1^\eta, \lambda)) := w(p_0^\eta) \cdot w(p_1^\eta)(n+1)^{-1}$ . Thus,  $\mathcal{P}'$  is the set of all two-component mixtures of  $\mathcal{P}^{(\eta)}$ ; and  $w'$  has its support on all two-component mixtures with  $\lambda \in \Lambda$ , and puts mass 0 on all other mixtures. Note that  $w'$  is indeed a prior, i.e.  $\sum_{p_0, p_1 \in \mathcal{P}, \lambda \in \Lambda} w'(\text{MIX}(p_0^\eta, p_1^\eta, \lambda)) \leq 1$ . We now apply Proposition 8 in the form (24) with, for all  $p' \in \mathcal{P}'$ ,  $\text{TR}(p') := q^\eta$ , and  $\eta$  in the proposition set to 1, and with the estimator that, for data  $z^n$ , chooses  $\text{MIX}(\hat{p}^\eta, \check{p}^\eta, \hat{\lambda})$ . That is, we set

$(\mathcal{P}, \check{p}, \mathcal{Q}, \text{TR}(\cdot), w, \eta) \leftarrow (\mathcal{P}', \text{MIX}(\hat{p}^\eta, \check{p}^\eta, \hat{\lambda}), \{q^\eta\}, \text{TR}(\cdot), w', 1)$ . This gives (41), where we also used (a), by definition,  $-\log w'(\text{MIX}(\hat{p}_0^\eta, \check{p}_1^\eta, \hat{\lambda})) \leq -\log w(\check{p})^2 + \log(n+1)$ ; and (b): since  $\eta \leq \eta_{\text{CRIT}}(u)$ , we have that  $d_1^*(q^\eta \| p^\eta) \geq -u/n$  for both  $p = \check{p}$  and  $p = \hat{p}$ , which implies, from the definition of  $d_1^*$ , that  $d_1^*(q^\eta \| (1-\lambda)\hat{p}^\eta + \lambda\check{p}^\eta) \geq -u/n$  for all  $\lambda \in [0, 1]$ , in particular for  $\hat{\lambda}$ .

Combining (40) and (41), using the definition of CONV-LACK, it follows that, WHP,

$$\begin{aligned} \text{CONV-LACK} &\leq c_\eta \text{RED}(2/\eta, \check{p}) + \frac{c_\eta}{\eta} (-\log w^2(\check{p}) + \log(n+1) + K + u + V^\eta) \\ &= c_\eta \text{RED}(4/\eta, \check{p}) + \frac{c_\eta}{\eta} (\log(n+1) + K + u + V^\eta). \end{aligned}$$

Now with some relatively straightforward manipulations\* we get that we get, WHP,  $\text{RED}(4/\eta, \check{p}) \leq 3\text{RED}(2/\eta, \check{p}) + \frac{2K+2u+2\log 2}{\eta}$ , so that, using the log-Bonferroni Proposition 12 with  $\mathcal{J} = 2$ , the above becomes  $\text{CONV-LACK} \leq 3c_\eta \text{RED}(2/\eta, \check{p}) + R$ , and the result follows.

**Proof of Lemma 5** A second-order Taylor expansion of  $E^*(p(Z_i)/q(Z_i))^\eta = E^*(e^{\eta \log p(Z_i)/q(Z_i)})$  around  $\eta = 0$  shows that, for all  $\eta > 0$ , for some  $0 \leq \eta' \leq \eta$ , we have:

$$E^* \left( \frac{p(Z_i)}{q(Z_i)} \right)^\eta = 1 - \eta E^*[-\log p(Z_i)/q(Z_i)] + \frac{1}{2} \eta^2 E^* \left( \frac{p(Z_i)}{q(Z_i)} \right)^{2\eta'} \left( \log \frac{p(Z_i)}{q(Z_i)} \right)^2 \leq 1 - \eta D^* + \frac{1}{2} \eta^2 V^{2\eta} D_{\text{SQ}}^*,$$

where we abbreviate  $D^*(q||p)$  to  $D^*$  and  $D_{\text{SQ}}^*(q, p)$  to  $D_{\text{SQ}}^*$ , and we replaced all factors in the expectation in the second order by their maximum. From now on we repeatedly use  $D^*(q||p) \geq 0$  which holds because  $q$  is best-approximating, It is sufficient to show that the right-hand side of this expression is bounded by  $1 + u/n$  if we plug in  $\eta \leq \eta_{\text{CRIT}}(u)$  as defined above. Dividing the inequality by  $\eta$  and using assumption (E), it is thus sufficient if we can show that

$$-D^* + \eta(D^*)^\nu \cdot b \leq \eta^{-1}(u/n) \quad (42)$$

where we set  $b = \frac{A}{2} V^{2\eta}$ . We may further assume  $(D^*)^{1-\nu} \leq \eta b$ , (43)

for if this does not hold, then  $-D^* = -(D^*)^\nu (D^*)^{1-\nu} \leq -(D^*)^\nu \eta b$  and then (42) holds trivially. Now first consider the case  $0 < \nu < 1$ . From (43) it follows that  $D^* \leq (\eta b)^{1/(1-\nu)}$ . By (42), it is thus sufficient if we can show that  $\eta \cdot (\eta b)^{\nu/(1-\nu)} b \leq \eta^{-1} u/n$ . Solving for  $\eta$  gives  $\eta^{2+\frac{\nu}{1-\nu}} \leq \frac{u}{n} b^{-1/(1-\nu)}$ , which can be rewritten to  $\eta \leq C$ , where  $C = \left(\frac{u}{n}\right)^{\frac{1-\nu}{2-\nu}} b^{-1/(2-\nu)}$ . Thus, weakening the requirement, it is sufficient if  $\eta \leq \min\{1/(2 \log V), C\}$ . But if  $\eta \leq 1/(2 \log V)$ , then  $b^{-1} \geq 2/(eA)$ , so it is also sufficient if  $\eta \leq \min\left\{\frac{1}{2 \log V}, B \left(\frac{u}{n}\right)^{\frac{1-\nu}{2-\nu}}\right\}$ . (20) now follows for the case  $0 < \nu < 1$ . The limiting cases  $\nu = 0$  and  $\nu = 1$  can be handled similarly; we omit details.

## 6. Discussion and Future Work

The great advances we made were already summarized on page 3; but currently, our work also has at least two major restrictions : (a)  $V = V(\mathcal{P}, P^*)$  as in (5) must be bounded; and (b)  $V$  occurs in the definition of CONV-LACK, so that it must be known in order to apply the safe estimator. Neither restriction is problematic for classification models, as long we used a fixed  $\beta$  in (2); both are problematic for e.g. standard regression models though. As to (a), currently our results only hold for such models if  $P^*$  has bounded support. In future work, we hope to replace the strong  $V < \infty$  condition with a weaker condition on moments of  $P^*$ . As to (b), we do have a version of all our results in which  $V$  is replaced by its empirical counterpart  $\bar{V} = \sup_{i \in \{1, \dots, n\}} \sup_{p, p' \in \mathcal{P}} p(Z_i)/p'(Z_i)$ , but with worse constants. We hope to refine this in future work.

Another issue is that, even if known,  $V$  or  $\bar{V}$ , appearing in CONV-LACK, may be so large as to make the approach useless in practice (even aside from computational issues, which, in this preliminary study, we decided not to deal with at all). We should note though that our current results hold for arbitrary priors  $w$ , in particular, priors with very heavy tails. Most priors used in practice have lighter tails, i.e.  $\sum_{p \in \mathcal{P}} w^\rho(p) < \infty$  for some  $\rho < 1$ . For such priors, the theorems still hold for the prior  $w'$  defined as  $w'(p) \propto w^\rho(p)$  rather

than the original  $w$ . As a result, the safe estimator  $\check{p}_{\text{SAFE}}$  defined relative to  $w'$  rather than  $w$  will effectively choose simpler distributions (with higher  $w(p)$ ) for the same data, but all occurrences of  $V$  in our theorems can be replaced by  $V^\rho$ , which can lead to a serious improvement in the size of CONV-LACK. A related idea is to consider the  $\beta$  in predictor models  $\mathcal{F}$  as in (2) as an additional parameter, to be equipped with a prior and fitted to the data. Since  $q$  and  $\check{p}$  in Theorem 1 may then refer to different predictors with different  $\beta$ 's,  $\beta$  will act as a ‘local’ learning rate whereas  $\eta$ , shared by all distributions, is a ‘global’ learning rate. Preliminary investigations suggest that this leads to better bounds in some cases.

## Acknowledgments

Supported in part by the IST Programme of the EU, under the PASCAL NoE, IST-2002-506778. I would like to thank Andrew Barron, Tim van Erven and Rui Castro for some very useful discussions.

## References

- J.Y. Audibert. *PAC-Bayesian statistical learning theory*. PhD thesis, Université Paris VI, 2004.
- A.R. Barron and T.M. Cover. Minimum complexity density estimation. *IEEE Transactions on Information Theory*, 37(4):1034–1054, 1991.
- A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth. Occam’s razor. *Information Processing Letters*, 24:377–380, 1987.
- L. Breiman. Statistical modeling: the two cultures (with discussion). *Statistical Science*, 16(3):199–215, 2001.
- O. Catoni. *PAC-Bayesian Supervised Classification*. Lecture Notes-Monograph Series. IMS, 2007.
- A.P. Dawid. The well-calibrated Bayesian. *Journal of the American Statistical Association*, 77:605–611, 1982. Discussion: pages 611–613.
- J.L. Doob. Application of the theory of martingales. In *Le Calcul de Probabilités et ses Applications. Colloques Internationaux du Centre National de la Recherche Scientifique*, pages 23–27, Paris, 1949.
- P. Grünwald. *The Minimum Description Length Principle*. MIT Press, Cambridge, MA, 2007.
- P. Grünwald and J. Langford. Suboptimal behavior of Bayes and MDL in classification under misspecification. *Machine Learning*, 66(2-3):119–149, 2007. DOI 10.1007/s10994-007-0716-7.
- P. D. Grünwald. Viewing all models as “probabilistic”. In *Proceedings of the Twelfth ACM Conference on Computational Learning Theory (COLT’ 99)*, pages 171–182, 1999.

- B. Kleijn and A. van der Vaart. Misspecification in infinite-dimensional Bayesian statistics. *Annals of Statistics*, 34(2), 2006.
- J.Q. Li. *Estimation of Mixture Models*. PhD thesis, Yale University, New Haven, CT, 1999.
- E. Mammen and A.B. Tsybakov. Smooth discrimination analysis. *Annals of Statistics*, 27: 1808–1829, 1999.
- D. McAllester. PAC-Bayesian stochastic model selection. *Machine Learning*, 51(1):5–21, 2003.
- A.B. Tsybakov. Optimal aggregation of classifiers in statistical learning. *Annals of Statistics*, 32:135–166, 2004.
- Tong Zhang. From  $\epsilon$ -entropy to KL entropy: analysis of minimum information complexity density estimation. *Annals of Statistics*, 34(5):2180–2210, 2006.

## Appendix A. Additional Proofs

**1. Condition (B) can be made to hold by adding one or a few distributions to  $\mathcal{P}$**  For example, in the classification case, it suffices to include the trivial distribution  $p_0$  into  $\mathcal{P}$ , where, for all  $x \in \mathcal{X}$ ,  $p_0(Y = 1 | X = x) = p_0(Y = 0 | X = x) = 1/2$ , and assign it some prior mass  $w_0(p_0) > 0$ . Then for all sequences  $z^n$ , for all  $\kappa \geq 1$ ,

$$\begin{aligned} -\log w(\tilde{p}_\kappa) &\leq -\log w(\tilde{p}_1) \leq -\log w(\tilde{p}_1) - \log \tilde{p}_1(z_1^n) \\ &\leq -\log w_0 - \log p_0(z_1^n) = -\log w_0 + n \log 2 \leq nL_{\max}, \end{aligned} \quad (44)$$

for suitably chosen  $L_{\max}$ . Clearly, this approach extends to all  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$  with finite or compact  $\mathcal{Y}$ . If  $\mathcal{Y}$  is not compact, then, by our assumption (D) that  $V(\mathcal{P}, P^*) < \infty$ , the interval  $[a, b]$  with  $a = \text{ess inf}_{Z \in \mathcal{Z}, p \in \mathcal{P}} p(Z_i)$  and  $b = \text{ess sup}_{p \in \mathcal{P}, Z \in \mathcal{Z}} p(Z_i)$ , is bounded. It then suffices to include a density  $p_{a,b}$  with prior  $w_0$  such that, for all  $x \in \mathcal{X}$ ,  $p_{a,b}(Y = \cdot | X = x)$  is uniform on  $[a, b]$ . If the end points on the interval are not known, we can discretize candidate end points to integers and put a prior  $v$  on these end points satisfying, for both end points  $c \in \{a, b\}$ ,  $-\log v(c) \approx 2 \log \min\{|c|, 1\}$  (Grünwald, 2007). We can define the defective distribution  $p_0(y | x) := \max_{a,b} p_{a,b}(y | x)v(a)v(b)$  and repeat the reasoning in (44).

**2. Equivalence of (4) and (12)** We will only show that equivalence holds in the idealized case in which the best-approximating  $q$  is actually a member of  $\mathcal{P}$ . This should be sufficient, since the goal of establishing equivalence is merely to give some intuition about the meaning of (12) (that’s why we only put it in the appendix); the equivalence is not needed in any of our results, whose proofs invariably rely on (12) rather than (4). Assume then that  $D(p^* || q) = \inf_{p \in \mathcal{P}} D(p^* || p)$  and that  $q \in \mathcal{P}$ . We only show equivalence for  $\eta = 1$ ; extension to other  $\eta$  is immediate. We first prove (4)  $\Rightarrow$  (12). If (4) holds, then for all  $p \in \mathcal{P}$ ,  $D^*(q || \text{MIX}(q, p, \lambda))$  has its minimum  $\lambda^\circ$  (as defined above Lemma 9) at  $\lambda^\circ = 0$ . It then follows by Lemma 9 that  $d_1^*(q || p) = d_1^*(q || \text{MIX}(q, p, 1)) \geq 0$ , and we see that (12) holds.

We next prove (12)  $\Rightarrow$  (4). Suppose that (12) holds for all  $p \in \mathcal{P}$ . Without loss of generality let  $\mathcal{P} = \{p_1, p_2, \dots\}$ . Then for any  $p'$  in the convex hull of  $\mathcal{P}$ , say  $p' = \sum_{j=1}^{\infty} \alpha_j p_j$

with all  $\alpha_j \geq 0$  and  $\sum \alpha_j = 1$ , we have  $E^*(p'(Z_i)/q(Z_i)) = \sum_{j=1}^{\infty} \alpha_j E^*(p_j(Z_i)/q(Z_i)) \leq 1$ . Thus  $E^*(p'(Z_i)/q(Z_i)) - 1 \leq 0$  and hence, by Proposition 10, part (1), the derivative of the concave function  $D^*(\text{MIX}(q, p', \lambda) \| q)$  is  $\leq 0$  at  $\lambda = 0$ . This implies that  $D(p^* \| q) \leq D(p^* \| \text{MIX}(q, p', \lambda))$  for all  $\lambda > 0$ , in particular  $D(p^* \| q) \leq D(p^* \| p')$ ; this shows that (4) holds.

**3. Lemma 9 – proof of Proposition 10** Differentiation gives:

$$g'(\gamma) = \frac{d}{d\gamma} g(\gamma) = -E^* \left( \frac{p(Z) - p'(Z)}{(1 - \gamma)p(Z) + \gamma p'(Z)} \right), \text{ in particular } g'(0) = E^* \left( \frac{p'(Z)}{p(Z)} \right) - 1, \quad (45)$$

which shows the first part of 1. We now first show part 2(a) and (b). Note that

$$g''(\gamma) = -E^* \left( \frac{p(Z) - p'(Z)}{(1 - \gamma)p(Z) + \gamma p'(Z)} \right)^2.$$

For all  $\gamma \in [0, 1]$  and all  $Z$ , the denominator inside the expectation must be bounded from below by  $\underline{p} := \text{ess inf}_{Z, p \in \mathcal{P}} p(Z)$  and from above by  $\bar{p} := \text{ess sup}_{Z, p \in \mathcal{P}} p(Z)$ . We thus have, for all  $\gamma \in [0, 1]$ ,  $1/\bar{p}^2 E^*(p' - p)^2 \leq |g''(\gamma)| \leq 1/\underline{p}^2 E^*(p' - p)^2$ . Now suppose first that  $E^*(p'(Z_i) - p(Z_i))^2 = 0$ . Then  $p'(Z_i) = p(Z_i)$  almost surely, and  $g'(\lambda) = g''(\lambda) = 0$  on  $[0, 1]$ , almost surely, and part (2a) follows. If  $E^*(p'(Z_i) - p(Z_i))^2 > 0$ , then  $p'(Z_i) \neq p(Z_i)$  with positive probability, and part (2b) immediately follows. Having now established that  $g''(\lambda) \leq 0$  on  $[0, 1]$ , it follows by definition of  $\lambda^\circ$  that  $g'(0) > 0$  iff  $\lambda^\circ > 0$ . And since we assume  $D^*(p \| p') \geq 0$ , we have  $g(1) \leq g(0)$ , which implies that if  $\lambda^\circ > 0$ , then  $g'(\lambda^\circ) = 0$  and  $g'(1) \leq 0$ .

**4. Proof of Proposition 12 (log-Bonferroni)** Let  $X_j := e^{-Y_j}$  and  $b_j = e^{-a_j}$ . The assumption implies that, for any collection  $\{K_j\}_{j \in \mathcal{J}}$  of positive real numbers, for all  $j \in \mathcal{J}$ ,

$$P^* \left( X_j \geq b_j e^{-f_j(K_j)} \right) \geq 1 - e^{-K_j},$$

or equivalently,

$$P^* \left( X_j < b_j e^{-f_j(K_j)} \right) < e^{-K_j}.$$

Now, for fixed  $K \geq 0$ , define  $K_j = K - \log w_j$ . By the union bound, we have

$$P^* (\mathcal{A}) < \sum_{j \in \mathcal{J}} e^{-K + \log w_j} = \sum_{j \in \mathcal{J}} w_j e^{-K}.$$

where  $\mathcal{A}$  is the event that for some  $j \in \mathcal{J}$ ,  $X_j < b_j e^{-f_j(K - \log w_j)}$ . This implies that for  $\bar{\mathcal{A}}$ , the complement of  $\mathcal{A}$ , we have

$$P^* (\bar{\mathcal{A}}) \geq 1 - \sum_{j \in \mathcal{J}} w_j e^{-K}.$$

The result now follows by noting that the event whose probability is bounded in the statement of the proposition is just  $\bar{\mathcal{A}}$ .

**5. Proof of Proposition 15 (Barron and Li's (1999) result)** Define, for given  $\eta, p^*, p$  and  $q$ , the *affinity relative to  $x$*  as  $A_x = \int_{y \in \mathcal{Y}} p^*(y | x) \cdot \left(\frac{p(y|x)}{q(y|x)}\right)^\eta$  and let

$$p^{\text{new}}(y | x) = \frac{1}{A(x)} p^*(y | x) \cdot \left(\frac{p(y | x)}{q(y | x)}\right)^\eta.$$

Next, recall that the *squared Hellinger distance*, between densities  $p$  and  $q$  on  $\mathcal{Y}$ , denoted by us as  $H^2(q, p)$ , is defined as

$$H^2(q, p) := \int_y (\sqrt{q(y)} - \sqrt{p(y)})^2 = 2 \left(1 - \int_y \sqrt{q(y)p(y)}\right).$$

Also recall that the ordinary (nongeneralized) Rényi divergence of order  $1/2$  is given by  $d_{1/2}(q, p) = -2 \log \int_y \sqrt{q(y)p(y)} dy$ . Now, for  $u \geq 0$ , we have  $1 - u \leq -\log u$  (this follows from  $\log(1+z) \geq z$  and substituting  $z = u - 1$ ). This implies the following well-known general relation between squared Hellinger distance and Rényi divergence:

$$H^2(q, p) \leq d_{1/2}(q \| p). \quad (46)$$

Moreover (Barron and Cover, 1991), when the ratio between  $p$  and  $q$  is bounded, then the standard (nongeneralized) KL divergence is upper-bounded by a multiple of the squared Hellinger distance. Yang and Barron (1999) proved the following precise relation:

$$D(q \| p) \leq (2 + \log V) H^2(q, p). \quad (47)$$

We will now use (46) and (47) to prove our result. We first need to clarify notation: for given  $x$ , the *generalized Rényi divergence between  $p$  and  $q$ , given  $x$*  is denoted as  $d_\eta^{*|x}(q(\cdot | x) \| p(\cdot | x))$  and defined as

$$d_\eta^{*|x}(q(\cdot | x) \| p(\cdot | x)) = -\frac{1}{\eta} \log E^* \left[ \left( \frac{p(Y | x)}{q(Y | x)} \right)^\eta \mid X = x \right].$$

We have for all  $C_\eta \geq 2 + 2\eta \log V$ , for each  $x \in \mathcal{X}$ , each  $\eta \leq 1$ ,

$$\begin{aligned} E^* \left[ -\log \frac{p(Y|x)}{q(Y|x)} \mid X = x \right] &= \frac{1}{\eta} \cdot E^* \left[ \log p^*(Y | x) - \log \left( p^*(Y | x) \left( \frac{q(Y|x)}{p(Y|x)} \right)^\eta \right) \mid X = x \right] \\ &\quad + \frac{1}{\eta} (\log A_x - \log A_x) \\ &= \frac{1}{\eta} D(p^*(\cdot|x) \| p^{\text{new}}(\cdot|x)) - \frac{1}{\eta} \log A_x \\ &\leq \frac{1}{\eta} C_\eta H^2(p^*(\cdot|x), p^{\text{new}}(\cdot|x)) - \frac{1}{\eta} \log A_x \\ &\leq \frac{1}{\eta} C_\eta d_{1/2}(p^*(\cdot|x), p^{\text{new}}(\cdot|x)) - \frac{1}{\eta} \log A_x \\ &= C_\eta \left( d_{\eta/2}^{*|x}(q(\cdot|x) \| p(\cdot|x)) + \frac{1}{\eta} \log A_x \right) - \frac{1}{\eta} \log A_x \\ &= C_\eta d_{\eta/2}^{*|x}(q(\cdot|x) \| p(\cdot|x)) + \frac{1}{\eta} (C_\eta - 1) \log A_x. \end{aligned}$$

Here the first two equalities are just rewriting. In the first inequality we used (47), the fact that  $P^*$ -almost surely,  $\sup_{X,Y} p^{\text{new}}(Y|X)/p^*(Y|X) \leq V^{2\eta}$ , and the fact that  $D(\cdot \| \cdot) \geq 0$ , and the second inequality is just (46). In the fifth line we used some basic rewriting. Using

the notation  $E_X^*$  to denote expectation of  $X$  under  $P_X^*$ , the marginal distribution of  $X$ , we thus get:

$$\begin{aligned}
 D^*(q\|p) &\leq C_\eta E_X^*[d_{\eta/2}^{*|X}(q(\cdot|X)\|p(\cdot|X))] + (C_\eta - 1)\frac{1}{\eta}E_X^*[\log A_X] \\
 &\leq C_\eta d_{2/\eta}^*(q\|p) + (C_\eta - 1)\frac{1}{\eta}\log E_X^*[A_X] \\
 &= C_\eta d_{\eta/2}^*(q\|p) - (C_\eta - 1)d_\eta^*(q\|p).
 \end{aligned}$$

where the second inequality is Jensen's and the final equality is just the definition of Rényi divergence.