

Mutual Information, Metric Entropy, and Risk in Estimation of Probability Distributions

*David Haussler**
UC Santa Cruz

Manfred Opper†
Universität Würzburg

December 29, 1996

University of California Technical Report UCSC-CRL-96-27
Baskin Center for Computer Science and Computer Engineering
UC Santa Cruz, CA 96064

This is a revision of the Dec. 27, 1995 version of this paper.
An abbreviated version of this paper is to appear in *Annals of Statistics*.

Abstract

Assume $\{P_\theta : \theta \in \Theta\}$ is a set of probability distributions with a common dominating measure on a complete separable metric space Y . A state $\theta^* \in \Theta$ is chosen by Nature. A statistician gets n independent observations Y_1, \dots, Y_n from Y distributed according to P_{θ^*} . For each time t between 1 and n , based on the observations Y_1, \dots, Y_{t-1} , the statistician produces an estimated distribution \hat{P}_t for P_{θ^*} , and suffers a loss $L(P_{\theta^*}, \hat{P}_t)$. The cumulative risk for the statistician is the average total loss up to time n . Of special interest in information theory, data compression, mathematical finance, computational learning theory and statistical mechanics is the special case when the loss $L(P_{\theta^*}, \hat{P}_t)$ is the relative entropy between the true distribution P_{θ^*} and the estimated distribution \hat{P}_t . Here the cumulative Bayes risk from time 1 to n is the mutual information between the random parameter Θ^* and the observations Y_1, \dots, Y_n .

New bounds on this mutual information are given in terms of the Laplace transform of the Hellinger distance between pairs of distributions indexed by parameters in Θ . From

*Supported by NSF grant IRI-9123692. Email addresses: haussler@cse.ucsc.edu

†Supported by Heisenberg fellowship of DFG. Email addresses: opper@physik.uni-wuerzburg.de

these, bounds on the cumulative minimax risk are given in terms of the metric entropy of Θ with respect to the Hellinger distance. The assumptions required for these bounds are very general and do not depend on the choice of the dominating measure. They apply to both finite and infinite dimensional Θ . They apply in some cases where Y is infinite dimensional, in some cases where Y is not compact, in some cases where the distributions are not smooth, and in some parametric cases where asymptotic normality of the posterior distribution fails. Using these bounds for cumulative relative entropy risk, we also examine the minimax risk of this game at specific times t for various loss functions L , including the relative entropy, the squared Hellinger distance, and the L_1 distance.

1 Introduction

Much of classical statistics has been concerned with the estimation of probability distributions from independent and identically distributed observations drawn according to these distributions. If we let P_{θ^*} denote the true distribution generating the observations and \hat{P}_t the estimated distribution obtained after seeing $t - 1$ independent observations, then the success of our statistical procedure can be defined in terms of a loss function that measures the difference between the true distribution P_{θ^*} and the estimated distribution \hat{P}_t . One such loss function has proven to be of importance in several fields, including information theory, data compression, mathematical finance, computational learning theory, and statistical mechanics. This is the relative entropy function. Further, in these fields, special importance is given to the cumulative relative entropy loss suffered in a sequential estimation setting, in which there are n total observations, but these observations arrive one at a time, and at each time t a new, refined estimate \hat{P}_t is made for the unknown true distribution P_{θ^*} , based on the $t - 1$ previous observations. This is the setting that we study in this paper.

The average of the cumulative loss over all sequences of n observations generated according to the true distribution is the (cumulative relative entropy) *risk*. For a given family $\{P_\theta : \theta \in \Theta\}$ of distributions, two types of risk are of interest in statistics. One is the minimax risk, which is the minimum worst-case risk over possible true distributions P_{θ^*} , where $\theta^* \in \Theta$, and the minimum is over all possible sequential estimation strategies. The other is the Bayes risk, which is the minimum average-case risk over possible true distributions P_{θ^*} drawn according to a prior distribution μ on Θ , and the minimum is again over all possible sequential estimation strategies. For cumulative relative entropy loss, the Bayes risk has a fundamental information theoretic interpretation: it is the mutual information between a random variable representing the choice of the parameter θ^* of the true distribution, and the random variable given by the n observations [37, 27, 18]. This provides a beautiful connection between information theory and statistics.

This connection also extends to other fields, as is discussed in [18, 8]. In data compression, the cumulative relative entropy risk is the *redundancy*, which is the expected excess code length for the best adaptive coding method, as compared to the best coding method that has prior knowledge of the true distribution [18, 41, 44]. The minimax risk is called “*information channel capacity*” [21], p. 184. In mathematical finance and gambling theory, the cumulative relative entropy risk measures the expected reduction in the logarithm of compounded wealth due to lack of knowledge of the true distribution [9, 18]. In computational learning theory,

this risk is the average additional loss suffered by an adaptive algorithm that predicts each observation before it arrives, based on the previous observations, as compared to an algorithm that makes predictions knowing the true distribution [34, 35]. Here we assume that the observation at time t is predicted by the “predictive” probability distribution \hat{P}_t , formed by the adaptive algorithm using the previous $t - 1$ observations, and that when this t th observation arrives, the loss is the negative logarithm of its probability under \hat{P}_t . Finally, in statistical mechanics, the Bayes risk can be related to the free energy [45, 46].

In this paper, we provide upper and lower bounds on the Bayes risk for cumulative relative entropy loss in the form of Laplace integrals of the Hellinger distance between pairs of distributions in $\{P_\theta : \theta \in \Theta\}$. We illustrate these bounds in a number of special cases, then use them to characterize the asymptotic rate of the minimax risk in terms of the metric entropy of $\{P_\theta : \theta \in \Theta\}$ under the Hellinger distance. The methods used here have the advantage of simplicity, with proofs amounting to little more than simple applications of Jensen’s inequality. The results are also quite general. The bounds apply to both finite and infinite dimensional Θ . They apply in some cases where the space of observations is infinite dimensional, in some cases where it is not compact, in some cases where the distributions are not smooth, and in some parametric cases where asymptotic normality of the posterior distribution fails. The bounds are also fairly tight. However, in smooth parametric cases, our general bounds are too crude to give the precise estimates of the low order additive constants that were obtained by Clarke and Barron [18, 19].

The paper is organized as follows. In sections 2 and 3 we give precise definitions of the risks that we evaluate, and discuss the conditions required for our bounds to hold. Here we also compare our bounds to those obtained previously by other authors. The bounds are given in section 4, followed by examples in sections 5 and 6 showing how they can be applied. Then in section 7 we give the characterization of the minimax risk. In sections 8 and 9 we illustrate further applications of our results by showing how they can be used to give bounds on the asymptotic rates for the minimax relative entropy risk at specific time t , as opposed to the cumulative risk. These results are then used further to obtain similar bounds for the risk under other loss functions, including the Hellinger and L_1 distance. Here the results are not as sharp as one can obtain by other methods, such as those of Le Cam [15, 42], Birgé [11, 12], Hasminskii and Ibragimov [32], and Wong and Shen [55], but these applications nevertheless illustrate the general utility of the method. Finally, we discuss some possible further work in section 10.

2 Basic definitions, notation and assumptions

The following notation and assumptions will be used throughout the paper.

Let Y be a complete separable metric space. All probability distributions on Y discussed in this paper are assumed to be defined on the σ -algebra of Borel sets of Y . Let Θ be a set, and for each $\theta \in \Theta$, let P_θ be a probability distribution on Y . We assume that for any $\theta \neq \theta^* \in \Theta$, the distributions associated with θ and θ^* are distinct in the sense that there is a Borel set $S \subset Y$ such that $P_\theta(S) \neq P_{\theta^*}(S)$. In addition, we assume there is a fixed σ -finite measure ν on Y that dominates P_θ for all $\theta \in \Theta$ (i.e. for any Borel set $S \subseteq Y$, $\nu(S) = 0$ implies $P_\theta(S) = 0$). We will also make (implicitly) the assumption that any other

distribution Q on Y mentioned in the results below is also dominated by ν . None of our results depend on the choice of the dominating measure ν , hence for any distribution Q , the Radon-Nikodym derivative $\frac{dQ}{d\nu}$ will be abbreviated simply as dQ , following the convention in Le Cam's text [42]. Furthermore, all integrals in the results below are assumed, without specific notation, to be taken with respect to the measure ν , unless otherwise indicated. Thus for a function f on Y and distribution Q on Y , the expectation of f is denoted

$$\int f dQ = \int_Y \frac{dQ}{d\nu}(y) f(y) d\nu.$$

Hence, in the special case that Y is countable and ν is the counting measure, for a probability mass function Q on Y

$$\int f dQ = \sum_{y \in Y} Q(y) f(y).$$

We will also need to treat probability distributions over Θ , which we will refer to as *prior distributions*. As each $\theta \in \Theta$ is associated with a distinct distribution P_θ on a complete separable metric space, we can define prior distributions on Θ with respect to the Borel sets of the topology of weak convergence of the P_θ measures. We assume that the set $\{P_\theta : \theta \in \Theta\}$ is itself measurable w.r.t. this topology. All prior distributions μ on Θ used in this paper are assumed to be Borel distributions of this type, and suprema over priors are also assumed to be only with respect to Borel distributions of this type. Further discussion of these issues can be found in the appendix of [25].

Finally, for integer or real-valued functions f and g , we say $f \sim g$ if $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 1$, and $f \asymp g$ if $\liminf_{n \rightarrow \infty} \frac{f(n)}{g(n)} > 0$ and $\limsup_{n \rightarrow \infty} \frac{f(n)}{g(n)} < \infty$. All logarithms are natural logarithms unless otherwise specified. We assume throughout that $0 \log 0 = 0 \log \frac{x}{0} = 0$, where x is any nonnegative finite number. We will also employ functions taking values in the extended reals $[-\infty, +\infty]$, and in particular use the extended log function obtained by defining $\log 0 = -\infty$ and $\log \infty = \infty$. Expectations over extended real-valued functions are defined whenever they do not take both the value $+\infty$ with positive probability and the value $-\infty$ with positive probability. The expectation is $+\infty$ if this value has positive probability, and similarly for $-\infty$.

3 Statement of the problem: the game of estimating a probability distribution

We view the problem of estimating a probability distribution from the set of distributions $\{P_\theta : \theta \in \Theta\}$ as a game in which Nature plays against the statistician. First Nature picks $\theta^* \in \Theta$. We refer to θ^* as the (*true*) *state of Nature*. Then for some $n \geq 1$, a sequence $Y^n = Y_1, \dots, Y_n$ of i.i.d. random variables are observed, each distributed according to P_{θ^*} . The particular sequence of values observed for these random variables is denoted $y^n = y_1, \dots, y_n$. For each time t between 1 and n , the statistician forms an estimate $\hat{P}_t = \hat{P}_t(y_t | y^{t-1})$ for the unknown distribution P_{θ^*} , based on the values $y^{t-1} = y_1, \dots, y_{t-1}$. In particular, for every t and every y^{t-1} , \hat{P}_t is a distribution over Y called the *predictive distribution at time t* , and

the set of all such predictive distributions, for all t and y^{t-1} , is called the (*predictive*) *strategy* of the statistician, and denoted simply as \hat{P} . Note that in this formulation, the statistician does not estimate the parameter θ^* itself, but rather the distribution it represents. This allows the statistician, if necessary, to use predictive distributions that are not in the set $\{P_\theta : \theta \in \Theta\}$.

Let L be a function that maps from pairs of distributions on Y into $[0, \infty]$. We call L the *loss function*. Specific loss functions we will consider include the *KL-divergence* or *relative entropy*, defined by

$$L(P, Q) = D_{KL}(P||Q) = \int dP \log \frac{dP}{dQ},$$

the (squared) Hellinger distance, defined by

$$L(P, Q) = D_{HL}^2(P, Q) = \int \left(\sqrt{dP} - \sqrt{dQ} \right)^2,$$

and the L_1 distance

$$L(P, Q) = \|P - Q\| = \sup_{|f| \leq 1} \left| \int f dP - \int f dQ \right| = \int |dP - dQ|.$$

For technical reasons, we will also consider the family of loss functions defined by the α -*affinities*, $\alpha > 1$, defined by

$$L(P, Q) = \rho_\alpha(P, Q) = \int (dP)^\alpha (dQ)^{1-\alpha}.$$

All of these loss functions are in the family of functions investigated by Csiszár, known as *f-divergences*, and all *f-divergences* are easily seen to be independent of the dominating measure [22]. The bulk of the paper is devoted to the relative entropy loss, so this loss is assumed unless otherwise specified.

For a fixed choice of loss function L , if the statistician uses the strategy \hat{P} , then the *risk (to the statistician) at time t , when θ^* is the state of Nature*, is given by

$$r_{t, \hat{P}, L}(\theta^*) = \int_{Y^{t-1}} dP_{\theta^*}^{t-1} L(P_{\theta^*}, \hat{P}_t).$$

The subscript L is omitted when the loss function is the relative entropy, here and in subsequent notation. The *cumulative risk for the first n observations* is

$$R_{n, \hat{P}, L}(\theta^*) = \sum_{t=1}^n r_{t, \hat{P}, L}(\theta^*).$$

The bulk of this paper discusses cumulative risk, which is henceforth referred to simply as *risk*, while the risk at time t is referred to as the *instantaneous risk*. For the relative entropy loss, the cumulative risk has a particularly simple interpretation. For any strategy \hat{P} , define the distribution \hat{P} on Y^n by

$$\hat{P}(y^n) = \prod_{t=1}^n \hat{P}_t(y_t | y^{t-1}).$$

In this way we can identify prediction strategies with joint distributions on Y_1, \dots, Y_n . Then

$$R_{n, \hat{P}}(\theta^*) = \sum_{t=1}^n \int_{Y^{t-1}} dP_{\theta^*}^{t-1}(y^{t-1}) \int_Y dP_{\theta^*}(y_t) \log \frac{dP_{\theta^*}(y_t)}{d\hat{P}_t(y_t|y^{t-1})} = D_{KL}(P_{\theta^*}^n || \hat{P}) \quad (1)$$

by the chain rule for relative entropy (see e.g. [21], p. 23).

Of course the statistician seeks a strategy that minimizes risk. One approach is to assume that Nature is a strategic adversary, and hence selects the worst case θ^* for any particular strategy of the statistician. In this case, the best strategy for the statistician is one that minimizes the worst-case risk, and the value of the game is the *minimax risk*

$$R_{n,L}^{minimax} = \inf_{\hat{P}} \sup_{\theta^* \in \Theta} R_{n, \hat{P}, L}(\theta^*).$$

A strategy \hat{P} that achieves this minimax value is called a *minimax strategy*. For the instantaneous risk, the corresponding minimax value is

$$r_{t,L}^{minimax} = \inf_{\hat{P}} \sup_{\theta^* \in \Theta} r_{t, \hat{P}, L}(\theta^*).$$

The other approach is the Bayesian approach, where one seeks to minimize the average risk. Here we might imagine that Nature chooses θ^* at random according to a prior probability distribution μ on Θ . Then the statistician seeks to minimize the average risk (according to μ), and the value of the game is the *Bayes risk*

$$R_{n, \mu, L}^{Bayes} = \inf_{\hat{P}} \int_{\Theta} d\mu(\theta^*) R_{n, \hat{P}, L}(\theta^*).$$

A strategy \hat{P} that achieves this value is called a *Bayes strategy*. For the instantaneous risk, the corresponding value is

$$r_{t, \mu, L}^{Bayes} = \inf_{\hat{P}} \int_{\Theta} d\mu(\theta^*) r_{t, \hat{P}, L}(\theta^*).$$

In the Bayesian approach there are two random variables, Θ^* , giving the choice of the state of Nature, and $Y^n = Y_1, \dots, Y_n$, giving the sequence of observations. Their joint distribution defines the behavior of Nature. The marginal distribution of Y^n , defined by

$$M_{n, \mu}(y^n) = \int_{\Theta} d\mu(\theta^*) P_{\theta^*}^n(y^n),$$

is of particular importance here. Breaking $M_{n, \mu}$ down into a product of conditional distributions, we can write

$$M_{n, \mu}(y^n) = \prod_{t=1}^n P_{t, \mu}^{Bayes}(y_t | y^{t-1}),$$

where

$$P_{t, \mu}^{Bayes}(y_t | y^{t-1}) = \frac{M_{t, \mu}(y^t)}{M_{t-1, \mu}(y^{t-1})}.$$

The distributions $P_{t,\mu}^{Bayes}$ are called *predictive posterior distributions*. These form a Bayes strategy for relative entropy loss, which we call P_μ^{Bayes} . To see this, note that by (1), the difference between the average cumulative risk for an arbitrary strategy \hat{P} and the strategy P_μ^{Bayes} is

$$\begin{aligned} \int_{\Theta} d\mu(\theta^*) \left(D_{KL}(P_{\theta^*}^n || \hat{P}) - D_{KL}(P_{\theta^*}^n || M_{n,\mu}) \right) &= \int_{\Theta} d\mu(\theta^*) \int_{Y^n} dP_{\theta^*}^n \left(\log \frac{dP_{\theta^*}^n}{d\hat{P}} - \log \frac{dP_{\theta^*}^n}{dM_{n,\mu}} \right) \\ &= \int_{Y^n} dM_{n,\mu} \log \frac{dM_{n,\mu}}{d\hat{P}} \\ &= D_{KL}(M_{n,\mu} || \hat{P}) \geq 0. \end{aligned}$$

It follows that the (cumulative) Bayes risk for relative entropy loss is given by

$$R_{n,\mu}^{Bayes} = \int_{\Theta} d\mu(\theta^*) D_{KL}(P_{\theta^*}^n || M_{n,\mu}) = I(\Theta^*; Y^n),$$

the *mutual information* between the parameter Θ^* and the observations Y^n . (See [21], p. 18, for general definition and discussion of the mutual information.)

It also turns out that for relative entropy loss, there is a simple, universal relationship between the Bayes risk $R_{n,\mu}^{Bayes}$ and the minimax risk $R_n^{minimax}$. This result can be obtained with limited effort from the general results in an early paper of Le Cam [14]. Special cases of the result were derived by Gallager [29] and Davisson and Leon-Garcia [23], and the general result is given in [33].

Theorem 1 [33]

$$R_n^{minimax} = \sup_{\mu} R_{n,\mu}^{Bayes},$$

where the supremum is taken over all (Borel) probability measures on the parameter space Θ . Moreover,

$$R_n^{minimax} = \inf_{\mu} \sup_{\theta^* \in \Theta} R_{n,P_\mu^{Bayes}}(\theta^*).$$

Several authors have studied the Bayes risk $R_{n,\mu}^{Bayes}$, or the equivalent mutual information $I(\Theta^*; Y^n)$, for the case of a parametric family of distributions $\{P_\theta : \theta \in \Theta\}$. Early work by Ibragimov and Hasminskii showed that $I(\Theta^*; Y^n) \sim (D/2) \log n$ when Y is the real line and the conditional distributions P_θ are a smooth family of densities indexed by a real-valued parameter vector θ in a compact set Θ of dimension D , and certain other conditions apply [37]. In this case they were even able to estimate the lower order additive terms in this approximation, which involve the Fisher information and the entropy of the prior. Further related results were given by Efroimovich [27] and Clarke [17]. Clarke and Barron gave a detailed analysis, with applications, of the risk of the Bayes strategy as a function of the true state of Nature [18], discussing the relation of the Bayes risk to the notion of redundancy in information theory, and giving applications to hypothesis testing and portfolio selection theory. These results were extended to the Bayes and minimax risk in [19] (see also [7]). Related lower bounds, which are often quoted, were obtained by Rissanen [51], based on certain asymptotic normality assumptions. Further extensions of this work are given by

Yamanishi [56, 58, 57]. Amari has developed an extensive theory that relates the risk when θ^* is the true state of Nature to certain differential-geometric properties of the parameter space Θ in the neighborhood of θ^* involving Fisher information and related quantities [2, 3] (see also [60, 40]).

Some authors have also looked at the value of the relative entropy risk in nonparametric cases as well, e.g. [6, 10, 52, 59, 55]. Also, the issue of consistent estimation of a general probability distribution with respect to relative entropy is addressed in [1, 41]. However, in the nonparametric case, more extensive work has been done in bounding the risk for other loss functions (see e.g. [24, 38]). While this work is too extensive to summarize here, we do note that some authors have also taken the general approach that we take here in using notions of metric entropy (defined below), and specifically using the Hellinger distance in obtaining these bounds (e.g. [42, 11, 12, 32, 54, 13, 10]). The only authors we have found who have applied this methodology to the relative entropy risk are Wong and Shen [55] (see Corollary 1, p. 360) and Barron and Yang [10]. This work is somewhat complementary to ours, in that it treats instantaneous risk, whereas we focus on cumulative risk. The tools that Wong and Shen employ are considerably more sophisticated, involving bracket entropy methods from empirical processes, and it appears that the boundedness assumptions they make (e.g. in Theorem 6) are a bit stronger than ours (see the discussion of integrable envelop functions at the end of section 4.2 below). Different assumptions, and different methods (using Fano's inequality) are used to obtain related general results in [10].

In this paper we describe a new approach, employing the Hellinger metric and certain Laplace integrals, to bounding both the Bayes and minimax risks for the cumulative relative entropy loss, and the instantaneous minimax risk for all three losses mentioned above.

The assumptions required for these general bounds are fairly mild. No special assumptions are needed for the lower bounds on the risk. To describe the assumptions needed for the upper bounds, recall that for $\alpha > 1$, $\rho_\alpha(P, Q) = \int (dP)^\alpha (dQ)^{1-\alpha}$. Hence, at time $t = 1$, i.e. when no observations have been made and the statistician must use some fixed *a priori* estimate \hat{P} for the true distribution P_{θ^*} , if the loss function is $L = \rho_{1+\lambda}$ for $\lambda > 0$ then the instantaneous minimax risk is the same as the cumulative risk for $n = 1$, and is given by

$$r_{1, \rho_{1+\lambda}}^{minimax} = R_{1, \rho_{1+\lambda}}^{minimax} = \inf_{\hat{P}} \sup_{\theta^*} \int (dP_{\theta^*})^{1+\lambda} (d\hat{P})^{-\lambda}.$$

For a prior distribution μ on Θ , the corresponding Bayes risk is

$$r_{1, \mu, \rho_{1+\lambda}}^{Bayes} = R_{1, \mu, \rho_{1+\lambda}}^{Bayes} = \inf_{\hat{P}} \int_{\Theta} d\mu(\theta^*) \int (dP_{\theta^*})^{1+\lambda} (d\hat{P})^{-\lambda}.$$

For our upper bounds on the minimax risk we make the assumption that there exists $\lambda > 0$ such that $R_{1, \rho_{1+\lambda}}^{minimax} < \infty$, and for the Bayes risk, that $R_{1, \mu, \rho_{1+\lambda}}^{Bayes} < \infty$. We also give an explicit formula for $R_{1, \mu, \rho_{1+\lambda}}^{Bayes}$. Further discussion of these assumptions, including some simple sufficient conditions for them to hold, and an artificial example in which they fail, is given at the end of section 4.2.

4 Bounds on mutual information and relative entropy distance to a mixture

Since we can obtain the minimax risk as a supremum of Bayes risks, we now focus our attention on the Bayes risk. As noted above, the Bayes risk $R_{n,\mu}^{Bayes}$ is the mutual information $I(\Theta^*, Y^n)$ between the random variable Θ^* giving the choice of θ^* according to the prior μ and the observations Y^n . We now give general bounds on this mutual information. In addition, since the risk for a particular state of Nature θ^* using the Bayes strategy P_μ^{Bayes} is

$$R_{n,P_\mu^{Bayes}}(\theta^*) = D_{KL}(P_{\theta^*}^n || M_{n,\mu}),$$

where $M_{n,\mu} = \int P_\theta^n d\mu(\theta)$, we will seek bounds for this quantity as well. The latter bounds actually address the general problem of bounding the relative entropy distance from an n -fold product distribution to a mixture of such distributions.

In obtaining these bounds, we use several notions of “distance” between probability distributions based on the α -affinities. One such family of distances are the I -divergences introduced by Renyi [49]. For any real $\alpha \neq 1$, and distributions P and Q , the I -divergence of order α is defined by

$$I_\alpha(P||Q) = \frac{1}{\alpha - 1} \log \int (dP)^\alpha (dQ)^{1-\alpha}. \quad (2)$$

For $0 < \alpha < 1$, a related set of distances is defined by

$$D_\alpha(P, Q) = \frac{1}{1 - \alpha} \left(1 - \int (dP)^\alpha (dQ)^{1-\alpha} \right) = \frac{1}{1 - \alpha} \int \left(\alpha dP + (1 - \alpha)dQ - (dP)^\alpha (dQ)^{1-\alpha} \right). \quad (3)$$

Since $\alpha x + (1 - \alpha)y - x^\alpha y^{1-\alpha} \geq 0$ for any $x, y \geq 0$ and $0 \leq \alpha \leq 1$, the integrand is everywhere nonnegative in the rightmost definition of D_α , showing that $D_\alpha(P, Q) \geq 0$. (This is essentially Hölder’s inequality.) Since $-\log x \geq 1 - x$, it follows that $I_\alpha(P||Q) \geq D_\alpha(P, Q)$, and hence $I_\alpha(P||Q) \geq 0$ as well. Since $-\log x \approx 1 - x$ for x near 1, these quantities are similar when the α -affinity $\int (dP)^\alpha (dQ)^{1-\alpha}$ is close to 1. Finally, for the case $\alpha = 1$, we define

$$D_1(P, Q) = I_1(P||Q) = D_{KL}(P||Q) = \int \left(dQ - dP - dP \log \frac{dQ}{dP} \right). \quad (4)$$

Since $\log z \leq z - 1$, it follows that $y - x - x \log \frac{y}{x} \geq 0$ for all $x, y \geq 0$, hence the integrand in the rightmost expression is everywhere nonnegative. It can be shown that both $D_\alpha(P, Q)$ and $I_\alpha(P||Q)$ are increasing in α for $\alpha > 0$.

One important special case of the above distances is the squared Hellinger distance

$$D_{HL}^2(P, Q) = D_{1/2}(P, Q) = \int (\sqrt{dP} - \sqrt{dQ})^2.$$

Unlike the other distances and divergences discussed above, the distance $D_{HL}(P, Q)$, i.e. the squareroot of the above defined D_{HL}^2 , is a metric, since it is symmetric and satisfies a triangle inequality. This metric has been used to give bounds on the risk of estimation procedures in statistics by many authors, including Le Cam [42], Birgé [11, 12], Hasminskii and Ibragimov [32], and van de Geer [54].

4.1 Basic bounds

Our main theorem gives bounds on $I(\Theta^*; Y^n)$ and $D_{KL}(P_{\theta^*}^n \| M_{n,\mu})$ in terms of the logarithms of two Laplace transforms of the I divergence, one at the value $\alpha = 1$ (the relative entropy) and the other at some α between 0 and 1.

Theorem 2 *Let μ be any prior measure on Θ and let $0 < \alpha < 1$. For each $\theta \in \Theta$ let Q_θ be an arbitrary conditional distribution on Y given θ and Q_θ^n be the n -fold product of Q_θ . For every $n \geq 1$,*

1.

$$\begin{aligned} - \int_{\Theta} d\mu(\theta^*) \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-n(1-\alpha)I_\alpha(P_{\theta^*} \| P_{\tilde{\theta}})} &\leq R_{n,\mu}^{Bayes} \\ &= I(\Theta^*; Y^n) \\ &\leq - \int_{\Theta} d\mu(\theta^*) \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-nI_1(P_{\theta^*} \| Q_{\tilde{\theta}})}. \end{aligned}$$

2. *For any $\gamma > 0$ there exists a subset Θ_γ of Θ with measure at least $1 - 2e^{-\gamma}$ under the prior μ such that for all $\theta^* \in \Theta_\gamma$*

$$\begin{aligned} - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-n(1-\alpha)I_\alpha(P_{\theta^*} \| P_{\tilde{\theta}})} - \gamma &\leq R_{n,P_\mu^{Bayes}}(\theta^*) \\ &= D_{KL}(P_{\theta^*}^n \| M_{n,\mu}) \\ &\leq - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-nI_1(P_{\theta^*} \| Q_{\tilde{\theta}})} + \gamma. \end{aligned}$$

The upper bound of part (1.) is similar to results given in [6], and is mentioned there for the case $P = Q$. To the best of our knowledge, the lower bound, and the results in part (2.), are new.

The proof is given in a series of lemmas and calculations. We prove the upper bounds of both parts of the theorem first, then the lower bounds. In establishing the bounds in part (2.), we will show that there is a set of μ -measure at most $e^{-\gamma}$ on which the lower bound fails, and similarly for the upper bound. Hence both bounds hold on the complement of the union of these two sets, which has μ -measure at least $1 - 2e^{-\gamma}$.

We begin with the upper bounds. This requires the following lemma which has been previously utilized in the framework of Statistical Physics [53].

Lemma 1 *Let $P = P(w)$ be a measure on a set W and $Q = Q(v)$ be a measure on a set V . For any real-valued function $u(w, v)$,*

$$- \int_V dQ(v) \log \int_W dP(w) e^{u(w,v)} \leq - \log \int_W dP(w) e^{\int_V dQ(v) u(w,v)}.$$

Proof: First note that by Hölder's inequality, for any real-valued functions u_1 and u_2 and $0 \leq \alpha \leq 1$,

$$\begin{aligned} \int_W dP(w) e^{\alpha u_1(w) + (1-\alpha)u_2(w)} &= \int_W dP(w) (e^{u_1(w)})^\alpha (e^{u_2(w)})^{(1-\alpha)} \\ &\leq \left(\int_W dP(w) e^{u_1(w)} \right)^\alpha \left(\int_W dP(w) e^{u_2(w)} \right)^{(1-\alpha)} \end{aligned}$$

Taking logs, this shows that $\log \int_W dP(w)e^{u(w,v)}$ is convex in u . The result then follows by applying Jensen's inequality. \square

We also use this simple lemma, suggested to us by Meir Feder. Let $P = P(v, w)$ be a measure on the product space $V \times W$, with conditional distribution $P(v|w)$ on V and marginal distribution $P(w)$ on W .

Lemma 2 *For any random variables W and V and nonnegative function $f(v, w)$ such that $\int_{V \times W} dP(v, w)f(v, w) = 1$,*

1.

$$\int_{V \times W} dP(v, w) \log f(v, w) \leq 0$$

2. For any $\gamma > 0$,

$$\Pr \left(w : \int_V dP(v|w) \log f(v, w) \geq \gamma \right) \leq e^{-\gamma}$$

Proof: For the first part, $\int_{V \times W} dP(v, w) \log f(v, w) = -\infty < 0$ if $f(v, w) = 0$ on a set of positive measure. Otherwise, note that by Jensen's inequality

$$\int_{V \times W} dP(v, w) \log f(v, w) \leq \log \int_{V \times W} dP(v, w)f(v, w) = 0.$$

Here we employ the convention that $0 \log 0 = 0$. For the second part, the case where $f(v, w) = 0$ for a set of v positive measure under the conditional distribution of V given w is similarly trivial, and otherwise note that

$$\begin{aligned} \Pr \left(w : \int_V dP(v|w) \log f(v, w) \geq \gamma \right) &= \Pr \left(w : e^{\int_V dP(v|w) \log f(v, w)} \geq e^\gamma \right) \\ &\leq e^{-\gamma} \int_W dP(w) e^{\int_V dP(v|w) \log f(v, w)} \\ &\leq e^{-\gamma} \int_W dP(w) \int_V dP(v|w) f(v, w) \\ &= e^{-\gamma} \end{aligned}$$

The first inequality follows from Markov's inequality and the second from Jensen's inequality. \square

In establishing the upper bounds, we use Lemma 2 with $V = Y^n$, $W = \Theta$ and $f(v, w) = \frac{\int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n(y^n)}{dM_{n, \mu}(y^n)}$. Here we assume all y^n such that $dM_{n, \mu} = 0$ have been removed from the domain of f , so that f is finite. The conditions of the lemma are satisfied, since this function is nonnegative and

$$\int_{\Theta \times Y^n} d\mu(\theta^*) dP_{\theta^*}^n(y^n) \frac{\int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n(y^n)}{dM_{n, \mu}(y^n)} = \int_{Y^n} \int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n(y^n) = 1,$$

since $M_{n, \mu}(y^n) = \int_{\Theta} d\mu(\theta^*) P_{\theta^*}^n(y^n)$. Employing Lemma 2 with this choice of f , the following chain of inequalities holds for all θ^* except for a set of μ -measure at most $e^{-\gamma}$.

$$D_{KL}(P_{\theta^*}^n || M_{n, \mu}) = \int_{Y^n} dP_{\theta^*}^n \log \frac{dP_{\theta^*}^n}{dM_{n, \mu}}$$

$$\begin{aligned}
&= \int_{Y^n} dP_{\theta^*}^n \left(\log \frac{dP_{\theta^*}^n}{\int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n} + \log \frac{\int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n}{dM_{n,\mu}} \right) \\
&\leq \int_{Y^n} dP_{\theta^*}^n \log \frac{dP_{\theta^*}^n}{\int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n} + \gamma \\
&= - \int_{Y^n} dP_{\theta^*}^n \log \int_{\Theta} d\mu(\tilde{\theta}) \frac{dQ_{\tilde{\theta}}^n}{dP_{\theta^*}^n} + \gamma \\
&= - \int_{Y^n} dP_{\theta^*}^n \log \int_{\Theta} d\mu(\tilde{\theta}) e^{\log \frac{dQ_{\tilde{\theta}}^n}{dP_{\theta^*}^n}} + \gamma \\
&\leq - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{\int_{Y^n} dP_{\theta^*}^n \log \frac{dQ_{\tilde{\theta}}^n}{dP_{\theta^*}^n}} + \gamma \\
&= - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-D_{KL}(P_{\theta^*}^n \| Q_{\tilde{\theta}}^n)} + \gamma \\
&= - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-nD_{KL}(P_{\theta^*} \| Q_{\tilde{\theta}})} + \gamma,
\end{aligned}$$

where the first inequality follows from Lemma (2) part (2.) and the second one from Lemma (1). The last equality follows from the fact that the KL divergence is additive over the product of independent distributions (see e.g. [21], p. 23). Note that by our convention that $0 \log 0 = 0$, for each θ^* , the set of y^n such that $dP_{\theta^*}^n(y^n) = 0$ can simply be removed in the first equality above and then reintroduced in the exponent of the second to the last inequality, thus avoiding any division by zero for these cases. Similarly, if $\int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n(y^n) = 0$ for a set of y^n of positive measure with respect to $P_{\theta^*}^n$, then all upper bounds from the second line on are infinite, and the result holds trivially. Otherwise a set of y^n of measure zero on which $\int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n(y^n) = 0$ can be ignored, avoiding any division by zero in this regard. Since $D_{KL} = I_1$, this establishes the upper bound of part (2.) of Theorem 2.

The upper bound of part (1.) of Theorem 2 is established in a very similar manner. Here we note that

$$\begin{aligned}
I(\Theta^*; Y^n) &= \int_{\Theta} d\mu(\theta^*) \int_{Y^n} dP_{\theta^*}^n \log \frac{dP_{\theta^*}^n}{dM_{n,\mu}} \\
&= \int_{\Theta} d\mu(\theta^*) \int_{Y^n} dP_{\theta^*}^n \left(\log \frac{dP_{\theta^*}^n}{\int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n} + \log \frac{\int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n}{dM_{n,\mu}} \right) \\
&\leq \int_{\Theta} d\mu(\theta^*) \int_{Y^n} dP_{\theta^*}^n \log \frac{dP_{\theta^*}^n}{\int_{\Theta} d\mu(\tilde{\theta}) dQ_{\tilde{\theta}}^n}
\end{aligned}$$

where the inequality follows from Lemma 2 part (1.). The remainder of the proof consists of the identical chain of inequalities as in the proof above of the upper bound of part (2.), except that we take expectation over θ^* and we do not have the term $+\gamma$.

We turn now to the lower bounds. Here we use the following lemma, which is new, as far as we can tell. Let $P = P(v, w)$ be a measure on the product space $V \times W$, with conditional

distribution $P(v|w)$ on V and marginal distribution $P(w)$ on W . For any $0 < \lambda \leq 1$, define¹

$$I^{(\lambda)}(W; V) = - \int_{V \times W} dP(v, w^*) \log \int_W dP(w) \left(\frac{dP(v|w)}{dP(v|w^*)} \right)^\lambda.$$

Note that $I^{(1)}(W; V) = I(W; V)$, the mutual information between W and V .

Lemma 3 *Whenever $\int_W dP(w)dP(v|w) > 0$ for all v , and $0 < \lambda \leq 1$,*

1.

$$I^{(\lambda)}(W; V) - I(W; V) \leq 0$$

2.

$$\Pr\{w^* : dP(v|w^*) > 0 \text{ and} \\ \int_V dP(v|w^*) \left(\log \int_W dP(w) \frac{dP(v|w)}{dP(v|w^*)} - \log \int_W dP(w) \left(\frac{dP(v|w)}{dP(v|w^*)} \right)^\lambda \right) \geq \gamma\} \leq e^{-\gamma}.$$

Proof: This follows from Lemma 2 using the function

$$f(v, w^*) = \frac{\int_W dP(w) \frac{dP(v|w)}{dP(v|w^*)}}{\int_W dP(w) \left(\frac{dP(v|w)}{dP(v|w^*)} \right)^\lambda} = \frac{dP^{\lambda-1}(v|w^*) \int_W dP(w) dP(v|w)}{\int_W dP(w) dP^\lambda(v|w)}.$$

The conditions of the lemma are satisfied, since f is nonnegative, $f(v, w^*) = 0$ when $dP(v|w^*) = 0$, and

$$\begin{aligned} \int_{V \times W} dP(v, w^*) f(v, w^*) &= \int_{V \times W} dP(v, w^*) \frac{dP^{\lambda-1}(v|w^*) \int_W dP(w) dP(v|w)}{\int_W dP(w) dP^\lambda(v|w)} \\ &= \int_V \int_W dP(w^*) dP(v|w^*) \frac{dP^{\lambda-1}(v|w^*) \int_W dP(w) dP(v|w)}{\int_W dP(w) dP^\lambda(v|w)} \\ &= \int_V \frac{\int_W dP(w^*) dP^\lambda(v|w^*) \int_W dP(w) dP(v|w)}{\int_W dP(w) dP^\lambda(v|w)} \\ &= \int_V \int_W dP(w) dP(v|w) \\ &= 1. \end{aligned} \tag{5}$$

¹For any v , $\int_W dP(w) \left(\frac{dP(v|w)}{dP(v|w^*)} \right)^\lambda = \frac{\int_W dP(w) dP^\lambda(v|w)}{dP^\lambda(v|w^*)} = 0$ or $= \frac{0}{0}$ only if $\int_W dP(w) dP(v|w) = 0$, which happens only if $dP(v|w) = 0$ for all but a set of w of measure zero. Using the convention that $-\log 0 = -\log \frac{0}{0} = 0$, the set of such v contribute nothing to $I^{(\lambda)}(W; V)$, and hence can be ignored. Furthermore, $\int_W dP(w) \left(\frac{dP(v|w)}{dP(v|w^*)} \right)^\lambda = \infty$ only if $P(v|w^*) = 0$ or $\int_W dP(w) dP^\lambda(v|w) = \infty$. However, by similar reasoning, for each individual w^* the set of all v such that $P(v|w^*) = 0$ can be ignored when evaluating $I^{(\lambda)}(W; V)$, and it is not possible that $\int_W dP(w) dP^\lambda(v|w) = \infty$, since $\int_W dP(w) dP(v|w) = dP(v) < \infty$ and for $0 < \lambda \leq 1$, $\int_W dP(w) dP^\lambda(v|w) \leq (\int_W dP(w) dP(v|w))^\lambda$ by Jensen's inequality. Thus $I^{(\lambda)}(W; V)$ is well-defined.

□

Now note that if $\{y^n : \int_{\Theta} d\mu(\tilde{\theta}) dP_{\tilde{\theta}}^n(y^n) = 0\}$ has positive measure under the distribution $P_{\tilde{\theta}^*}^n$ then $D_{KL}(P_{\tilde{\theta}^*}^n || M_{n,\mu}) = \infty$. Hence the lower bound holds trivially. Otherwise a set of such y^n of measure zero can be ignored, and using part (2.) of Lemma 3 with $W = \Theta$ and $V = Y^n$, we can show that the following inequalities hold except on a set of θ^* with μ -measure at most $e^{-\gamma}$.

$$\begin{aligned}
D_{KL}(P_{\tilde{\theta}^*}^n || M_{n,\mu}) &= - \int_{Y^n} dP_{\tilde{\theta}^*}^n \log \int_{\Theta} d\mu(\tilde{\theta}) \frac{dP_{\tilde{\theta}}^n}{dP_{\tilde{\theta}^*}^n} \\
&\geq - \int_{Y^n} dP_{\tilde{\theta}^*}^n \log \int_{\Theta} d\mu(\tilde{\theta}) \left(\frac{dP_{\tilde{\theta}}^n}{dP_{\tilde{\theta}^*}^n} \right)^{\lambda} - \gamma \\
&\geq - \log \int_{\Theta} d\mu(\tilde{\theta}) \int_{Y^n} dP_{\tilde{\theta}^*}^n \left(\frac{dP_{\tilde{\theta}}^n}{dP_{\tilde{\theta}^*}^n} \right)^{\lambda} - \gamma \\
&= - \log \int_{\Theta} d\mu(\tilde{\theta}) \int_{Y^n} (dP_{\tilde{\theta}^*}^n)^{1-\lambda} (dP_{\tilde{\theta}}^n)^{\lambda} - \gamma \\
&= - \log \int_{\Theta} d\mu(\tilde{\theta}) \left[\int_Y (dP_{\tilde{\theta}^*})^{1-\lambda} (dP_{\tilde{\theta}})^{\lambda} \right]^n - \gamma \\
&= - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-n \log \int_Y (dP_{\tilde{\theta}^*})^{1-\lambda} (dP_{\tilde{\theta}})^{\lambda}} - \gamma \\
&= - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-n \lambda I_{1-\lambda}(P_{\tilde{\theta}} || P_{\tilde{\theta}^*})} - \gamma.
\end{aligned}$$

As in the proof of the upper bound, to avoid division by zero, and apply Lemma 3, we can remove the set of y^n such that $dP_{\tilde{\theta}^*}^n(y^n) = 0$ from the first line, and reintroduce them in the fourth line. Setting $\alpha = 1 - \lambda$, this establishes the lower bound of part (2.).

As with the upper bound, the lower bound of part (1.) is established easily by removing the $-\gamma$ terms and taking expectation over θ^* in the above chain of inequalities, using part (1.) of Lemma 3 in line 2. This establishes the lower bounds, and completes the proof of Theorem 2. □

A few brief comments about Theorem 2 are in order. First, note that if in part (2) we let γ grow with n in a suitable way, we obtain bounds which asymptotically hold for almost all $\theta^* \in \Theta$. An even stronger result is obtained when we chose $\gamma(n)$ such that $\sum_{n=1}^{\infty} e^{-\gamma(n)} < \infty$. This holds for example, if we let $\gamma(n)$ grow faster than $\log n$. Then, the first Borel-Cantelli lemma shows that for μ almost all $\theta \in \Theta$, the bounds will be violated only a finite number of times as $n \rightarrow \infty$.

It should also be noted that in the important special case when $P = Q$, the upper bound of part (2) of the theorem holds with $\gamma = 0$, since we can omit the first few steps of its derivation in this case, where γ is introduced. Thus both this strengthened upper bound and the given lower bound hold on a set of measure $1 - e^{-\gamma}$ in this case.

Finally, we note that part (2) is related to part (1) in the same way that the strong redundancy-capacity theorem of universal coding in [44] is related to the usual theorems concerning average redundancy.

It is possible to state a variant of Theorem 2 using the the D_{α} distances. Here we also make use of a particular choice for the family of distributions Q_{θ} that appear in Theorem 2.

Another possible choice is explored in Theorem 3 below. We will need the following definition.

For each $0 < \alpha < 1$ and $x > 0$ define

$$b_\alpha(x) = \frac{(1-\alpha)(x - \log x - 1)}{\alpha + (1-\alpha)x - x^{1-\alpha}}. \quad (6)$$

Define $b_\alpha(0) = \infty$. It is easily verified that $b_\alpha(x)$ is strictly decreasing in x , approaches 1 as $x \rightarrow \infty$, and approaches ∞ as $x \rightarrow 0$. Let

$$B_\alpha(\Theta) = \sup_{y \in Y, \theta^*, \tilde{\theta} \in \Theta} b_\alpha \left(\frac{dP_{\theta^*}(y)}{dP_{\tilde{\theta}}(y)} \right).$$

Clearly this constant does not depend on the choice of the dominating measure ν .

Corollary 1 *For every $0 < \alpha < 1$ and $n \geq 1$,*

1.

$$\begin{aligned} - \int_{\Theta} d\mu(\theta^*) \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-n(1-\alpha)D_\alpha(P_{\theta^*}, P_{\tilde{\theta}})} &\leq R_{n,\mu}^{Bayes} \\ &= I(\Theta^*; Y^n) \\ &\leq - \int_{\Theta} d\mu(\theta^*) \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-nB_\alpha(\Theta)D_\alpha(P_{\theta^*}, P_{\tilde{\theta}})}. \end{aligned}$$

2. *For any $\gamma > 0$ there exists a subset Θ_γ of Θ with measure at least $1 - 2e^{-\gamma}$ under the prior μ such that for all $\theta^* \in \Theta_\gamma$*

$$\begin{aligned} - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-n(1-\alpha)D_\alpha(P_{\theta^*}, P_{\tilde{\theta}})} - \gamma &\leq R_{n,P_\mu^{Bayes}}(\theta^*) \\ &= D_{KL}(P_{\theta^*}^n \| M_{n,\mu}) \\ &\leq - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-nB_\alpha(\Theta)D_\alpha(P_{\theta^*}, P_{\tilde{\theta}})} + \gamma. \end{aligned}$$

Proof. Since $I_\alpha(P||Q) \geq D_\alpha(P, Q)$, the lower bounds follow directly from the lower bounds of Theorem 2. For the upper bounds, we will need the following lemma, which is a simple extension of Lemma 4.4 of [11].

Lemma 4 *For any distributions P and Q on Y and any $0 < \alpha < 1$,*

$$D_{KL}(P||Q) \leq \left(\sup_{y \in Y} b_\alpha \left(\frac{dQ(y)}{dP(y)} \right) \right) D_\alpha(P, Q).$$

Proof. If $dP = dQ$ except on a set of zero measure (w.r.t. the dominating measure ν), then $D_{KL}(P||Q) = 0$, and hence the result holds. So it suffices to consider the case where $D_\alpha(P, Q) > 0$. Let $S = \{y \in Y : dP(y) = 0\}$. Separating Y into S and $Y - S$, and factoring a dP out of the integrands in Equations (3) and (4) in the latter case, we have

$$\begin{aligned} \frac{D_{KL}(P||Q)}{D_\alpha(P, Q)} &= \frac{(1-\alpha) \int_{Y-S} dP \left(\frac{dQ}{dP} - \log \frac{dQ}{dP} - 1 \right) + \int_S dQ}{\int_{Y-S} dP \left(\alpha + (1-\alpha) \frac{dQ}{dP} - \left(\frac{dQ}{dP} \right)^{1-\alpha} \right) + \int_S dQ} \\ &\leq \sup_{y \in Y} b_\alpha \left(\frac{dQ(y)}{dP(y)} \right). \end{aligned}$$

since $b_\alpha \geq 1$. \square

The upper bounds of Corollary 1 follow from Theorem 2 and this lemma by setting $Q_\theta = P_\theta$. \square

Whenever dP_θ is uniformly bounded above zero and below infinity for all y and θ for some choice of the dominating measure, $B_\alpha(\Theta)$ is finite, and this corollary can be applied. However, in some other cases $B_\alpha(\Theta) = \infty$ for all $0 < \alpha < 1$, making the upper bound in the above corollary useless. One case where this occurs is when there are θ and θ^* in Θ such that P_θ is not dominated by P_{θ^*} . For example, if $Y = \{0, 1\}$ and there is a θ^* such that $P_{\theta^*}(Y = 1)$ is zero (or one) and there is also a θ where $P_\theta(Y = 1)$ is not zero (or not one), then P_θ is not dominated by P_{θ^*} . We can also have $B_\alpha(\Theta) = \infty$ in cases where such lack of domination does not occur. For example, if $Y = \{0, 1\}$, Θ is the open interval $(0, 1)$, and $P_\theta(Y = 1) = \theta$, then $B_\alpha(\Theta) = \infty$ not because there are two distributions that fail to mutually dominate each other, but because $\inf_{y \in Y, \theta^*, \theta \in \Theta} \frac{dP_{\theta^*}(y)}{dP_\theta(y)} = 0$. Such cases can be handled by the results in the following section.

4.2 Bounds for finite $R_{1,\mu,\rho_{1+\lambda}}^{Bayes}$

Here we prove a version of Corollary 1 that can be used in cases when $B_\alpha(\Theta) = \infty$ for all $0 < \alpha < 1$. This new theorem requires only the weaker assumption that the Bayes risk for the $(1 + \lambda)$ -affinity loss at time 1, $R_{1,\mu,\rho_{1+\lambda}}^{Bayes}$, discussed in section 3 above, is finite for some $\lambda > 0$. Recall that for a fixed prior μ ,

$$R_{1,\mu,\rho_{1+\lambda}}^{Bayes} = \inf_{\hat{P}} \int_{\Theta} d\mu(\theta^*) \int (dP_{\theta^*})^{1+\lambda} (d\hat{P})^{-\lambda}.$$

Using Jensen's inequality, it can be verified that when $R_{1,\mu,\rho_{1+\lambda}}^{Bayes} < \infty$, the minimizing \hat{P} , i.e. the Bayes strategy, is the distribution $U = U_\mu$ defined by

$$dU = \frac{\left(\int_{\Theta} d\mu(\theta^*) dP_{\theta^*}^{1+\lambda} \right)^{\frac{1}{1+\lambda}}}{C_{\lambda,\mu}},$$

where

$$C_{\lambda,\mu} = \int_Y \left(\int_{\Theta} d\mu(\theta^*) dP_{\theta^*}^{1+\lambda} \right)^{\frac{1}{1+\lambda}}$$

[60]. Hence for each individual θ^* , the risk of the Bayes strategy is

$$\begin{aligned} R_{1,U_\mu,\rho_{1+\lambda}}(\theta^*) &= \int_Y (dP_{\theta^*})^{1+\lambda} (dU)^{-\lambda} \\ &= C_{\lambda,\mu}^\lambda \int_Y (dP_{\theta^*})^{1+\lambda} \left(\int_{\Theta} d\mu(\theta^*) dP_{\theta^*}^{1+\lambda} \right)^{-\frac{\lambda}{1+\lambda}} \end{aligned}$$

and the Bayes risk is

$$\begin{aligned} R_{1,\mu,\rho_{1+\lambda}}^{Bayes} &= \int_{\Theta} d\mu(\theta^*) \int_Y (dP_{\theta^*})^{1+\lambda} (dU)^{-\lambda} \\ &= C_{\lambda,\mu}^\lambda \int_{\Theta} d\mu(\theta^*) \int_Y (dP_{\theta^*})^{1+\lambda} \left(\int_{\Theta} d\mu(\theta^*) dP_{\theta^*}^{1+\lambda} \right)^{-\frac{\lambda}{1+\lambda}} \end{aligned}$$

$$\begin{aligned}
&= C_{\lambda,\mu}^\lambda \int_Y \left(\int_{\Theta} d\mu(\theta^*) dP_{\theta^*}^{1+\lambda} \right)^{\frac{1}{1+\lambda}} \\
&= C_{\lambda,\mu}^{1+\lambda}
\end{aligned}$$

We have the following theorem.

Theorem 3 *Let $0 < \alpha < 1$ and $0 < \lambda \leq 1$. Assume $R_{1,\mu,\rho_{1+\lambda}}^{Bayes} < \infty$. Then for every $n \geq 1$,*

1.

$$\begin{aligned}
& - \int_{\Theta} d\mu(\theta^*) \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-n(1-\alpha)D_\alpha(P_{\theta^*}, P_{\tilde{\theta}})} \\
& \leq R_{n,\mu}^{Bayes} \\
& = I(\Theta^*; Y^n) \\
& \leq - \int_{\Theta} d\mu(\theta^*) \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-(n \log n) \frac{(1+o(1))4(1-\alpha)}{\alpha\lambda} D_\alpha(P_{\theta^*}, P_{\tilde{\theta}})} + R_{1,\mu,\rho_{1+\lambda}}^{Bayes} + o(1)
\end{aligned}$$

2. *For any $\gamma > 0$ there exists a subset Θ_γ of Θ with measure at least $1 - 2e^{-\gamma}$ under the prior μ such that for all $\theta^* \in \Theta_\gamma$*

$$\begin{aligned}
& - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-n(1-\alpha)D_\alpha(P_{\theta^*}, P_{\tilde{\theta}})} - \gamma \\
& \leq R_{n,P_\mu^{Bayes}}(\theta^*) \\
& = D_{KL}(P_{\theta^*}^n \| M_{n,\mu}) \\
& \leq - \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-(n \log n) \frac{(1+o(1))4(1-\alpha)}{\alpha\lambda} D_\alpha(P_{\theta^*}, P_{\tilde{\theta}})} + R_{1,U,\rho_{1+\lambda}}(\theta^*) + \gamma + o(1),
\end{aligned}$$

where in each case for fixed α and λ , $o(1)$ is a function $f(n)$ such that $f(n) \rightarrow 0$ and $n \rightarrow \infty$. Furthermore, the same results also hold replacing the quantity $D_\alpha(P_{\theta^*}, P_{\tilde{\theta}})$ with $I_\alpha(P_{\theta^*} \| P_{\tilde{\theta}})$.

Proof. That D can be replaced by I follows from the fact that $I_\alpha(P \| Q) \geq D_\alpha(P, Q)$ for all α , as pointed out in the proof of Corollary 1. To prove the result for D , we will need a lemma².

Lemma 5 *Assume $0 < \alpha < 1$ and $\lambda > 0$. Let P , R and U be any distributions on Y . Let $c_\lambda = \int dP^{1+\lambda} dU^{-\lambda}$. Let $Q = (1 - \epsilon)R + \epsilon U$ for some $\epsilon > 0$ such that $\frac{\log \log(1/\epsilon)}{\log(1/\epsilon)} \leq \lambda/2$ and $\epsilon \leq e^{-\alpha/(2(1-\alpha))}$. Then*

$$D_{KL}(P \| Q) \leq \frac{2 \log \frac{1}{\epsilon}}{f_\alpha(\epsilon^2)} D_\alpha(P, R) + \frac{2\epsilon \log \frac{1}{\epsilon}}{(1-\alpha)f_\alpha(\epsilon^2)} + \epsilon^{\lambda/2} c_\lambda,$$

where

$$f_\alpha(x) = \frac{\alpha + (1-\alpha)x - x^{1-\alpha}}{1-\alpha}.$$

²We recently noticed that a related result is given in [55], Theorem 5, although no explicit relationship with the α affinities is given in the latter result.

The proof of this lemma is given in the appendix.

Now let U_μ be the Bayes strategy as defined above. Since $R_{1,\mu,\rho_{1+\lambda}}^{Bayes} < \infty$, U_μ is well-defined. For each $\theta \in \Theta$, let

$$Q_\theta = (1 - \epsilon)P_\theta + \epsilon U_\mu,$$

with $\epsilon = n^{-2/\lambda}$. It is clear that $f_\alpha(\epsilon^2) \rightarrow \frac{\alpha}{1-\alpha}$ as $\epsilon \rightarrow 0$. Hence, by Lemma 5, for sufficiently large n , for all θ

$$\begin{aligned} D_{KL}(P_{\theta^*} || Q_{\hat{\theta}}) &\leq \frac{2 \log \frac{1}{\epsilon}}{f_\alpha(\epsilon^2)} D_\alpha(P_{\theta^*}, P_{\hat{\theta}}) + \frac{2\epsilon \log \frac{1}{\epsilon}}{(1-\alpha)f_\alpha(\epsilon^2)} + \epsilon^{\lambda/2} R_{1,U_\mu,\rho_{1+\lambda}}(\theta^*) \\ &= \frac{4 \log n}{\lambda f_\alpha(n^{-4/\lambda})} D_\alpha(P_{\theta^*}, P_{\hat{\theta}}) + \frac{R_{1,U_\mu,\rho_{1+\lambda}}(\theta^*) + o(1)}{n} \quad \text{since } \lambda \leq 1 \\ &= \log n \frac{(1 + o(1))4(1-\alpha)}{\alpha \lambda} D_\alpha(P_{\theta^*}, P_{\hat{\theta}}) + \frac{R_{1,U_\mu,\rho_{1+\lambda}}(\theta^*) + o(1)}{n} \end{aligned}$$

Since $\int_\Theta d\mu(\theta^*) R_{1,U_\mu,\rho_{1+\lambda}}(\theta^*) = R_{1,\mu,\rho_{1+\lambda}}^{Bayes}$, the result then follows from Theorem 2. \square

Note: no attempt has been made to optimize the constants in this theorem.

Now let³

$$S(\Theta) = \int_Y \sup_{\theta \in \Theta} dP_\theta.$$

We call $\sup_{\theta \in \Theta} dP_\theta$ the *envelop function* for Θ . Note that $S(\Theta)$ is independent of the choice of the dominating measure. Since for all $\lambda \geq 0$,

$$\left(\int_\Theta d\mu(\theta^*) dP_{\theta^*}^{1+\lambda} \right)^{\frac{1}{1+\lambda}} \leq \sup_{\theta^* \in \Theta} dP_{\theta^*},$$

It follows that

$$R_{1,\mu,\rho_{1+\lambda}}^{Bayes} = C_{\lambda,\mu}^{1+\lambda} = \left(\int_Y \left(\int_\Theta d\mu(\theta^*) dP_{\theta^*}^{1+\lambda} \right)^{\frac{1}{1+\lambda}} \right)^{1+\lambda} \leq S^{1+\lambda}(\Theta)$$

for all $\lambda > 0$. Hence, whenever Θ has an integrable envelop function, that is whenever $S(\Theta) < \infty$, then $R_{1,\mu,\rho_{1+\lambda}}^{Bayes} < \infty$, and the bounds in part (1) of Theorem 3 hold with $\lambda = 1$ and $R_{1,\mu,\rho_{1+\lambda}}^{Bayes}$ replaced with $S^2(\Theta)$. It is clear that $S(\Theta) < \infty$ whenever Y is finite, and whenever Y is a bounded set in R^k for some $k \geq 1$ and the densities in $\{P_\theta : \theta \in \Theta\}$ are uniformly upper bounded. Hence Theorem 3 always applies in these cases.

Theorem 3 also applies in many cases where $S(\Theta)$ is infinite; an example of such a case is given in the following section. To characterize the types of Θ and priors ν not covered by Theorem 3, let us define the function $f_{\Theta,\mu}(\lambda)$ for $\lambda \geq 0$ by

$$f_{\Theta,\mu}(\lambda) = \frac{1}{\lambda} \log R_{1,\mu,\rho_{1+\lambda}}^{Bayes}$$

³If $\sup_{\theta \in \Theta} dP_\theta$ is not measurable, then any measurable function that majorizes it can be used instead in the definition of $S(\Theta)$.

for $\lambda > 0$ and

$$f_{\Theta, \mu}(0) = R_{1, \mu}^{Bayes},$$

that is, the risk at time 1 for the relative entropy loss. It can be shown that for any Θ and μ , $f_{\Theta, \mu}(\lambda)$ is a nondecreasing function on $[0, \infty)$ taking values in $[0, \infty]$, and if $f_{\Theta, \mu}(\lambda)$ is finite for any $\lambda > 0$, then

$$\lim_{\lambda \rightarrow 0} f_{\Theta, \mu}(\lambda) = f_{\Theta, \mu}(0).$$

To verify this last property, note that $\lim_{\lambda \rightarrow 0} R_{1, \mu, \rho_{1+\lambda}}^{Bayes} = 1$. Hence by l'hospital's rule

$$\lim_{\lambda \rightarrow 0} f_{\Theta, \mu}(\lambda) = \frac{d}{d\lambda} \left(R_{1, \mu, \rho_{1+\lambda}}^{Bayes} \right) \Big|_{\lambda=0}.$$

It can be verified by direct calculation that the latter quantity is the mutual information $I(\Theta^*, Y)$, which is the same as $R_{1, \mu}^{Bayes}$.

It is clear that whenever $R_{1, \mu}^{Bayes}$ is infinite, then $R_{n, \mu}^{Bayes}$ is infinite for all $n \geq 1$. Thus there are only three possible cases for the pair (Θ, μ) :

1. $f_{\Theta, \mu}(\lambda) < \infty$ for some $\lambda > 0$. In this case $R_{1, \mu, \rho_{1+\lambda}}^{Bayes} < \infty$ and hence Theorem 3 applies and may be used to get bounds on $R_{n, \mu}^{Bayes}$ for all n .
2. $f_{\Theta, \mu}(0) = \infty$. In this case $R_{n, \mu}^{Bayes} = \infty$ for all n and hence the problem of bounding this quantity is trivial.
3. $f_{\Theta, \mu}(0) < \infty$ but $f_{\Theta, \mu}(\lambda) = \infty$ for all $\lambda > 0$. In this case we say that the pair (Θ, μ) is *irregular*. These are the only nontrivial cases where Theorem 3 does not apply.

While it would not be expected that irregular (Θ, μ) would show up much in practice, it is possible to construct one.

Example 1 Let $Y = \{1, 2, 3, \dots\}$, $\Theta = \{3, 4, 5, \dots\}$, and for each $\theta \in \Theta$ and $y \in Y$, define $P_\theta(Y = y)$ to be $1 - \frac{1}{\log \theta}$ if $y = 1$, $\frac{1}{\log \theta}$ if $y = \theta$, and 0 otherwise. Let $\mu(\theta) = \frac{c}{\theta \log^2 \theta}$, where $c = \sum_{i=3}^{\infty} \frac{1}{i \log^2 i} < \infty$. Then it can be shown that (Θ, μ) is irregular.

5 Examples

We now illustrate Theorems 2 and 3 by applying them to a few simple problems. We begin with a classical case in which each point $\theta \in \Theta$ is a vector of D real numbers, Θ is a compact set and the prior μ is specified as a density $d\mu(\theta)$. To apply Theorem 2, fix $\theta^* \in \Theta_\gamma$, where θ^* is in the interior of Θ . We assume that the prior $d\mu$ is continuous and positive at θ^* . We also assume that $\{P_\theta\}$ is a smooth family of probabilities such that the *Fisher information* matrix at θ^* , defined by $J(\theta^*)$, where

$$J_{ij}(\theta^*) = \int_Y dP_{\theta^*} \left[\frac{\partial}{\partial \theta_i} \log dP_\theta(y) \frac{\partial}{\partial \theta_j} \log dP_\theta(y) \right] \Big|_{\theta=\theta^*},$$

exists and is positive definite. In this case, we will focus on the bounds on the risk for individual θ^* , rather than bounds on the mutual information. Even the simplest choice $Q = P$

will be sufficient to obtain a useful bound in the smooth case. For large n , obviously the main contributions to the inner expectations in Theorem 2 come from small neighborhoods of θ^* . Hence, under certain regularity conditions, Laplace's method can be used to evaluate these expectations asymptotically. We perform a Taylor expansion of the exponents in Theorem 2 to second order in the difference between $\tilde{\theta}$ and θ^* using the partial derivatives

$$\frac{\partial}{\partial \theta_i} I_\alpha(P_{\theta^*} \| P_\theta)|_{\theta=\theta^*} = 0$$

and

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} I_\alpha(P_{\theta^*} \| P_\theta)|_{\theta=\theta^*} = \alpha J_{ij}(\theta^*). \quad (7)$$

Note, that these results are also valid for $\alpha = 1$. Hence, Laplace's method would yield for the lower bound

$$\int_{\Theta} d\mu(\tilde{\theta}) e^{-n(1-\alpha)I_\alpha(P_{\theta^*} \| P_{\tilde{\theta}})} = d\mu(\theta^*) \int_{\mathbb{R}^D} d\theta e^{-\frac{n}{2}\alpha(1-\alpha)\sum_{ij}(\theta_i-\theta_i^*)J_{ij}(\theta^*)(\theta_i-\theta_i^*)}(1+o(1)).$$

A similar expression is obtained for the upper bound. By evaluating the Gaussian integrals we get⁴

$$\begin{aligned} \frac{D}{2} \log \frac{n}{2\pi} - \log d\mu(\theta^*) + \frac{1}{2} \log \det J(\theta^*) - \frac{D}{2} \log \frac{1}{\alpha(1-\alpha)} - \gamma + o(1) &\leq R_{n, P_\mu^{\text{Bayes}}}(\theta^*) \leq \\ &\frac{D}{2} \log \frac{n}{2\pi} - \log d\mu(\theta^*) + \frac{1}{2} \log \det J(\theta^*) + o(1). \end{aligned}$$

Note that asymptotically the lower bound is optimized by setting $\alpha = \frac{1}{2}$. In this case, for large n , both bounds differ by a constant approximately equal to $\frac{D \log 4}{2}$ for small γ . In this classical case, Clarke and Barron [18] have determined the exact answer to within $o(1)$, and it is

$$R_{n, P_\mu^{\text{Bayes}}}(\theta^*) = \frac{D}{2} \log \frac{n}{2\pi} - \log d\mu(\theta^*) + \frac{1}{2} \log \det J(\theta^*) - \frac{D}{2} + o(1).$$

Thus our simpler methods do not give the best known additive constants in the bounds for this classical case, but they do provide good bounds for large n .

As pointed out by Clarke and Barron [18], the scaling $\sim \frac{D}{2} \log n$ of the Bayes risk for the smooth parametric families is strongly related to the asymptotic normality of the properly normalized posterior distribution. It is interesting to look at nonregular families of probabilities, for which the posterior fails to converge to a nontrivial limit. (For conditions that are necessary for convergence, see [30]). As an example for such nonsmooth densities, we study the following simple family on \mathbb{R}

$$dP_\theta(y) = e^{-(y-\theta)} I_{\{y>\theta\}}, \quad \theta \in \mathbb{R}. \quad (8)$$

Obviously, $D_{KL}(P_{\theta^*} \| P_\theta) = \infty$, whenever $\theta > \theta^*$ and the Fisher information does not exist for any θ . Hence, the previous analysis is not applicable and we have to resort to the more sophisticated upper bounds. Specializing to $\alpha = \frac{1}{2}$, we easily find

$$\begin{aligned} D_{1/2}(P_{\theta^*}, P_\theta) &= 2(1 - e^{-|\theta-\theta^*|}) \\ I_{1/2}(P_{\theta^*} \| P_\theta) &= |\theta - \theta^*|. \end{aligned}$$

⁴Here we can set $\gamma = 0$ in the upper bounds, as per the comments following Theorem 2.

This result clearly shows the difference from the smooth families. The distances $D_{1/2}$ and $I_{1/2}$ do not behave locally like a quadratic function for θ close to θ^* , but have a linear scaling. Hence, a different scaling of the risk at θ^* and the mutual information is also expected.

An explicit result using Theorem 3 is easily obtained for the prior $d\mu(\theta) = \frac{1}{2}e^{-|\theta|}$. Note that the envelope of Θ is not integrable, so we must obtain direct bounds on $R_{1,\mu,\rho_{1+\lambda}}$ rather than using $S(\Theta)$. To upper bound $R_{1,\mu,\rho_{1+\lambda}}^{Bayes} = \inf_{\hat{P}} \int_{\Theta} d\mu(\theta^*) f(dP_{\theta^*})^{1+\lambda} (d\hat{P})^{-\lambda}$ it suffices to choose any distribution U and bound the expectation of $c_{\lambda}(\theta^*) = f(dP_{\theta^*})^{1+\lambda} (dU)^{-\lambda}$. Here we can set $dU(y) = \frac{1}{2}e^{-|y|}$. In this case we have $c_{\lambda}(\theta^*) < e^{\lambda|\theta^*|}$ and $\int_{\Theta} d\mu(\theta^*) c_{\lambda}(\theta^*) < \infty$ for all $\lambda < 1$. To evaluate the bounds we use the fact that for $a > 1$

$$\frac{1}{2} \int_{-\infty}^{\infty} d\theta e^{-|\theta|-a|\theta-\theta^*|} = \frac{e^{-|\theta^*|} - e^{-a|\theta^*|}}{2(a-1)} + \frac{e^{-|\theta^*|} + e^{-a|\theta^*|}}{2(a+1)}.$$

Hence, for $\alpha = \frac{1}{2}$, we get

$$\log\left(\frac{n}{2}\right) + |\theta^*| - \gamma + o(1) \leq R_{n,P_{\mu}^{Bayes}}(\theta^*) \leq \log\left(\frac{4n \log n(1+o(1))}{\lambda}\right) + e^{\lambda|\theta^*|} + |\theta^*| + \gamma + o(1).$$

Hence, an asymptotic scaling $\sim \log n$ for the risk is observed. This gives a factor of two difference compared to the risk of a smooth 1-dimensional family of densities.

Finally, we will consider an example where both the parameter space and the space of observations are infinite dimensional. We assume that an unknown real continuous function $\theta(x)$ with $0 \leq x \leq 1$ is corrupted by a Gaussian white noise process. The statistician observes n random functions Y_t , $t = 1, \dots, n$ which, conditioned on θ , are independent realizations of the process

$$Y(x) = \int_0^x \theta(z) dz + \sigma W(x). \quad (9)$$

Here $W(x)$ is a standard Wiener process with $W(0) = 0$ and covariance $\mathbb{E}[W(x_1)W(x_2)] = \min(x_1, x_2)$. In this case, it is easy to calculate the I-divergences explicitly for all α . Let P_{θ} be the measure corresponding to the random process $Y(x)$ and let the dominating measure ν be the Wiener measure. Then, from the Cameron–Martin formula [16], the Radon-Nikodym derivative is found to be

$$\frac{dP_{\theta}}{d\nu} = \exp\left[\frac{1}{\sigma} \int_0^1 \theta(x) dW(x) - \frac{1}{2\sigma^2} \int_0^1 \theta^2(x) dx\right]. \quad (10)$$

Inserting this into the definition of the I-divergences, we obtain

$$I_{\alpha}(P_{\theta^*} || P_{\theta}) = \frac{\alpha}{2\sigma^2} \int_0^1 (\theta(x) - \theta^*(x))^2 dx. \quad (11)$$

For the case where the prior over the space of functions $\theta(x)$ is a Gaussian measure (such that $\theta(x)$ is a realization of a Gaussian random process) our bounds can be evaluated in closed form. We will restrict ourselves to the case of the mutual information $I(\Theta^*; Y^n)$ and use the fact that for Gaussian processes and $c > 0$

$$-\int_{\Theta} d\mu(\theta^*) \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-\frac{c}{2} \int_0^1 (\tilde{\theta}(x) - \theta^*(x))^2 dx} = \frac{1}{2} \sum_k \left[\log(1 + c\lambda_k) + \frac{c\lambda_k}{1 + c\lambda_k} \right]. \quad (12)$$

Here $\lambda_k, k = 1, 2, \dots, \infty$ are the eigenvalues of the process on the interval $[0, 1]$. Specializing on the Wiener process, we get $\lambda_k = \frac{1}{\pi^2(k-\frac{1}{2})^2}$, for $k = 1, 2, 3 \dots$. Using

$$\frac{1}{2} \sum_{k=1}^{\infty} \log\left(1 + \frac{c}{\pi^2(k-\frac{1}{2})^2}\right) = \frac{1}{2} \log \cosh(\sqrt{c})$$

and

$$\frac{1}{2} \sum_{k=1}^{\infty} \frac{c}{c + \pi^2(k-\frac{1}{2})^2} = \frac{\sqrt{c}}{4} \tanh \sqrt{c}$$

and setting $\alpha = \frac{1}{2}$, we get

$$\begin{aligned} \frac{1}{2} \log \cosh\left(\frac{\sqrt{n}}{2\sigma}\right) + \frac{\sqrt{n}}{8\sigma} \tanh\left(\frac{\sqrt{n}}{2\sigma}\right) &\leq I(\Theta^*; Y^n) \leq \\ \frac{1}{2} \log \cosh\left(\frac{\sqrt{n}}{\sigma}\right) + \frac{\sqrt{n}}{4\sigma} \tanh\left(\frac{\sqrt{n}}{\sigma}\right). & \end{aligned}$$

Hence, asymptotically

$$\frac{3\sqrt{n}}{8}(1 + o(1)) \leq I(\Theta^*; Y^n) \leq \frac{3\sqrt{n}}{4}(1 + o(1)).$$

Notice that in the above examples, it was always the case that asymptotically, the best bounds were obtained with the value $\alpha = 1/2$. In general, for large n the value of the Laplace transform

$$\int_{\Theta} d\mu(\tilde{\theta}) e^{-n(1-\alpha)I_{\alpha}(P_{\theta^*}||P_{\tilde{\theta}})}$$

in the lower bound of Theorem 2 is largely determined by those $\tilde{\theta}$ such that $I_{\alpha}(P_{\theta^*}||P_{\tilde{\theta}})$ is near zero, i.e. such that $P_{\tilde{\theta}}$ is close to P_{θ^*} . The same also holds for the corresponding Laplace transform

$$\int_{\Theta} d\mu(\tilde{\theta}) e^{-nI_1(P_{\theta^*}||P_{\tilde{\theta}})}$$

in the upper bound. However, it can be shown that as the distributions P and Q get close, in the sense that $\frac{dP}{dQ} \rightarrow 1$ uniformly, then

$$\frac{I_1(P||Q)}{(1-\alpha)I_{\alpha}(P||Q)} \rightarrow \frac{1}{\alpha(1-\alpha)}.$$

Hence we might expect to very often get the best asymptotic lower bound in Theorem 2 by choosing $\alpha = 1/2$, so as to minimize $\frac{1}{\alpha(1-\alpha)}$. This choice also has another desirable property, since, as mentioned above, for $\alpha = 1/2$, the distance D_{α} used in Corollary 1 and Theorem 3 is then the squared Hellinger distance, which has some nice metric properties that we will exploit in applications of the bounds below. For these reasons, in what follows, we will for simplicity restrict ourselves to the case $\alpha = 1/2$, using the notation

$$D_{1/2}(P, Q) = D_{HL}^2(P, Q).$$

6 Bounds on the cumulative risk for countable Θ

Recall that we have assumed that for all distinct $\theta, \theta^* \in \Theta$, the conditional densities dP_θ and dP_{θ^*} differ on a set of positive measure, and hence $D_{HL}(P_\theta, P_{\theta^*}) > 0$. We can make this assumption without essential loss of generality, since otherwise we can replace Θ by a set of equivalence classes with the property that $\theta \equiv \theta^*$ iff $dP_\theta = dP_{\theta^*}$ (except on a set of measure zero) in a natural way, without changing the risks we are interested in calculating.

Suppose Θ is countable, say $\Theta = \{\theta_i\}$. Let $H(\Theta^*) = -\sum_i \mu(\theta_i) \log \mu(\theta_i)$ denote the entropy of the random variable Θ^* , distributed according to the prior measure μ . The entropy of Θ^* may be infinite. Then

Corollary 2 *For all n , $R_{n,\mu}^{Bayes} = I(\Theta; Y^n) \leq H(\Theta^*)$ and*

$$\lim_{n \rightarrow \infty} R_{n,\mu}^{Bayes} = H(\Theta^*).$$

Proof: Recall that $R_{n,\mu}^{Bayes} = I(\Theta^*; Y^n)$. If $H(\Theta^*)$ is infinite then clearly

$$\limsup_{n \rightarrow \infty} I(\Theta; Y^n) \leq H(\Theta^*).$$

Assume $H(\Theta^*)$ is finite. Let

$$H(\Theta^* | Y^n) = - \int_{Y^n} dM_{n,\mu}(y^n) \sum_i \mu(\theta_i | y^n) \log \mu(\theta_i | y^n),$$

the conditional entropy of Θ given Y^n . Note that this quantity is nonnegative. When $H(\Theta)$ is finite it is easily verified that

$$I(\Theta^*; Y^n) = H(\Theta^*) - H(\Theta^* | Y^n)$$

(see e.g. [21], p. 20), and thus $\limsup_{n \rightarrow \infty} I(\Theta^*; Y^n) \leq H(\Theta^*)$ in this case as well.

For the lower bound, using Theorem 2 with $\alpha = 1/2$ and Fatou's lemma

$$\begin{aligned} \liminf_{n \rightarrow \infty} I(\Theta^*; Y^n) &\geq \liminf_{n \rightarrow \infty} - \sum_i \mu(\theta_i) \log \sum_j \mu(\theta_j) e^{-\frac{n}{2} D_{HL}^2(P_{\theta_i}, P_{\theta_j})} \\ &\geq - \sum_i \mu(\theta_i) \liminf_{n \rightarrow \infty} \log \sum_j \mu(\theta_j) e^{-\frac{n}{2} D_{HL}^2(P_{\theta_i}, P_{\theta_j})} \\ &= - \sum_i \mu(\theta_i) \log \mu(\theta_i) \\ &= H(\Theta^*). \end{aligned}$$

□

This result generalizes the similar result in [19] (Corollary 1) by removing the additional conditions assumed there. More general results, including the above corollary, follow from results in Pinsker's book [47] (see also [4]). Applying Theorem (1) and taking the supremum over μ in Corollary (2), it follows that if Θ is finite then for all n , $R_n^{minimax} \leq \log |\Theta|$ and $\lim_{n \rightarrow \infty} R_n^{minimax} = \log |\Theta|$. It also follows that if Θ is infinite, then $\lim_{n \rightarrow \infty} R_n^{minimax} = \infty$.

In the case that Θ is finite, results of Renyi [50] show further that the difference $I(\Theta^*; Y^n) - H(\Theta^*)$ converges to zero exponentially fast in n . We also obtain this result as follows.

Corollary 3 For all n ,

$$H(\Theta^*) - I(\Theta^*; Y^n) \leq (|\Theta| - 1) \left(\max_{1 \leq i < j \leq |\Theta|} \int_Y \sqrt{dP_{\theta_i} dP_{\theta_j}} \right)^n.$$

Proof: From Theorem 2

$$\begin{aligned} I(\Theta^*; Y^n) &\geq - \sum_i \mu(\theta_i) \log \sum_j \mu(\theta_j) \left(\int_Y \sqrt{dP_{\theta_i} dP_{\theta_j}} \right)^n \\ &= - \sum_i \mu(\theta_i) \log \mu(\theta_i) - \sum_i \mu(\theta_i) \log \left[1 + \sum_{j \neq i} \frac{\mu(\theta_j)}{\mu(\theta_i)} \left(\int_Y \sqrt{dP_{\theta_i} dP_{\theta_j}} \right)^n \right] \\ &\geq H(\Theta^*) - \sum_i \sum_{j \neq i} \mu(\theta_j) \left(\int_Y \sqrt{dP_{\theta_i} dP_{\theta_j}} \right)^n \\ &\geq H(\Theta^*) - (|\Theta| - 1) \left(\max_{1 \leq i < j \leq |\Theta|} \int_Y \sqrt{dP_{\theta_i} dP_{\theta_j}} \right)^n, \end{aligned}$$

where the second inequality follows from $-\log(1+x) \geq -x$. \square

Assuming as above that the densities dP_{θ_i} and dP_{θ_j} are different for $j \neq i$, an application of the Cauchy's inequality yields $\int_Y \sqrt{dP_{\theta_i} dP_{\theta_j}} < 1$ for $j \neq i$. Hence, the corollary shows exponential convergence.

Finally, let us note that Theorem 2 and Corollary 1 can also be used to characterize the mutual information between Θ^* and Y^n (Bayes risk) in the general case when Θ is uncountably infinite but finite dimensional. This was demonstrated in [36]. Here in the sequel, we focus instead on the minimax risk.

7 Bounds on minimax risk using covering and packing numbers, and metric entropy

For each $\theta^*, \theta \in \Theta$, let

$$h(\theta^*, \theta) = D_{HL}(P_{\theta^*}, P_{\theta}).$$

As mentioned above, we assume that for distinct states of Nature $\theta, \theta^* \in \Theta$, the conditional distributions P_{θ} and P_{θ^*} differ on a set of positive measure. Under this assumption, (Θ, h) is a metric space. We show how bounds on the minimax risk can be obtained by looking at properties of this metric space. These are the the packing and covering numbers, and the associated metric entropy, introduced by Kolmogorov and Tikhomirov in [39] and commonly used in the theory of empirical processes (see e.g. [26, 48, 31, 13]).

For the following definitions, let (S, ρ) be any complete separable metric space.

Definition 1 (*Metric entropy, also called Kolmogorov ϵ -entropy [39]*) A partition Π of S is a collection $\{\pi_i\}$ of Borel subsets of S that are pairwise disjoint and whose union is S . The diameter of a set $A \subseteq S$ is given by $\text{diam}(A) = \sup_{x, y \in A} \rho(x, y)$. The diameter of a partition is the supremum of the diameters of the sets in the partition. For $\epsilon > 0$, by $\mathcal{D}_{\epsilon}(S, \rho)$ we

denote the cardinality of the smallest finite partition of S of diameter at most ϵ , or ∞ if no such finite partition exists. The metric entropy of (S, ρ) is defined by

$$\mathcal{K}_\epsilon(S, \rho) = \log \mathcal{D}_\epsilon(S, \rho).$$

We say S is totally bounded if $\mathcal{D}_\epsilon(S, \rho) < \infty$ for all $\epsilon > 0$.

Definition 2 (*Packing and covering numbers*) For $\epsilon > 0$, an ϵ -cover of S is a subset $A \subseteq S$ such that for all $x \in S$ there exists a $y \in A$ with $\rho(x, y) \leq \epsilon$. By $\mathcal{N}_\epsilon(S, \rho)$ we denote the cardinality of the smallest finite ϵ -cover of S , or ∞ if no such finite cover exists. For $\epsilon > 0$, an ϵ -separated subset of S is a subset $A \subseteq S$ such that for all distinct $x, y \in A$, $\rho(x, y) > \epsilon$. By $\mathcal{M}_\epsilon(S, \rho)$ we denote the cardinality of the largest finite ϵ -separated subset of S , or ∞ if arbitrarily large such sets exist.

The following lemma is easily verified [39].

Lemma 6 For any $\epsilon > 0$,

$$\mathcal{M}_{2\epsilon}(S, \rho) \leq \mathcal{D}_{2\epsilon}(S, \rho) \leq \mathcal{N}_\epsilon(S, \rho) \leq \mathcal{M}_\epsilon(S, \rho).$$

It follows that the metric entropy \mathcal{K}_ϵ (and the condition defining total boundedness) can also be defined using either the packing or covering numbers in place of \mathcal{D}_ϵ , to within a constant factor in ϵ .

Kolmogorov and Tikhomirov also introduced abstract notions of the dimension and order of metric spaces in their seminal paper [39]. These can be used to measure the ‘‘massiveness’’ of both spaces indexed by a finite dimensional parameter vector and infinite dimensional function spaces. In the following, the metric ρ is omitted from the notation, being understood from the context.

Definition 3 The upper and lower metric dimensions [39] of S are defined by

$$\overline{\mathbf{dim}}(S) = \limsup_{\epsilon \rightarrow 0} \frac{\mathcal{K}_\epsilon(S)}{\log \frac{1}{\epsilon}}$$

and

$$\underline{\mathbf{dim}}(S) = \liminf_{\epsilon \rightarrow 0} \frac{\mathcal{K}_\epsilon(S)}{\log \frac{1}{\epsilon}},$$

respectively. When $\overline{\mathbf{dim}}(S) = \underline{\mathbf{dim}}(S)$, then this value is denoted $\mathbf{dim}(S)$ and called the metric dimension of S . Thus

$$\mathbf{dim}(S) = \lim_{\epsilon \rightarrow 0} \frac{\mathcal{K}_\epsilon(S)}{\log \frac{1}{\epsilon}}.$$

For totally bounded S , we say that S is finite dimensional if $\mathbf{dim}(S) < \infty$, else it is infinite dimensional. To measure the massiveness of infinite dimensional spaces, including typical function spaces, further indices were introduced by Kolmogorov and Tikhomirov. The functional dimension of S is defined similarly as

$$\mathbf{df}(S) = \lim_{\epsilon \rightarrow 0} \frac{\log \mathcal{K}_\epsilon(S)}{\log \log \frac{1}{\epsilon}},$$

with similar upper and lower versions when this limit does not exist. Finally, the metric order of S is defined as

$$\mathbf{mo}(S) = \lim_{\epsilon \rightarrow 0} \frac{\log \mathcal{K}_\epsilon(S)}{\log \frac{1}{\epsilon}},$$

with similar upper and lower versions.

Using the results given in the theorems from section 4, with $\alpha = 1/2$, we can obtain bounds on the minimax risk R_n^{minimax} in terms of the metric entropy of the space (Θ, h) . For every $\epsilon > 0$ let

$$b(\epsilon) = \sup \left\{ \frac{D_{KL}(P_{\tilde{\theta}} \| P_{\theta^*})}{D_{HL}^2(P_{\tilde{\theta}}, P_{\theta^*})} : \tilde{\theta}, \theta^* \in \Theta \text{ and } D_{HL}^2(P_{\tilde{\theta}}, P_{\theta^*}) \leq \epsilon \right\}.$$

Recall also that the minimax risk for time 1 and loss $\rho_{1+\lambda}$ is denoted $R_{1, \rho_{1+\lambda}}^{\text{minimax}}$.

Lemma 7 *Assume (Θ, h) is totally bounded. Then for all $n \geq 1$,*

1.

$$\begin{aligned} R_n^{\text{minimax}} &\geq \sup_{\epsilon \geq 0} \left\{ -\log \left(\frac{1}{\mathcal{M}_\epsilon(\Theta, h)} + e^{-\frac{n\epsilon^2}{2}} \right) \right\} \\ &\geq \sup_{\epsilon \geq 0} \min \left\{ \mathcal{K}_\epsilon(\Theta, h), \frac{n\epsilon^2}{8} \right\} - \log 2 \end{aligned}$$

and

2.

$$R_n^{\text{minimax}} \leq \inf_{\epsilon \geq 0} \left\{ \mathcal{K}_\epsilon(\Theta, h) + b(\epsilon)n\epsilon^2 \right\} \leq \inf_{\epsilon \geq 0} \left\{ \mathcal{K}_\epsilon(\Theta, h) + b_{1/2}(\Theta)n\epsilon^2 \right\}.$$

Furthermore, for any $\lambda > 0$ such that $R_{1, \rho_{1+\lambda}}^{\text{minimax}} < \infty$

$$R_n^{\text{minimax}} \leq \inf_{\epsilon \geq 0} \left\{ \mathcal{K}_\epsilon(\Theta, h) + \frac{(1 + o(1))4\epsilon^2 n \log n}{\lambda} \right\} + R_{1, \rho_{1+\lambda}}^{\text{minimax}} + o(1),$$

where in each case $o(1)$ is a function $f(n)$ such that $f(n) \rightarrow 0$ as $n \rightarrow \infty$.

Proof: To establish the first inequality of part (1), let $A = \{\theta_1, \dots, \theta_M\}$ be an ϵ -separated subset of Θ of maximal size and let μ be the discrete prior distribution on Θ that is uniform over the elements of A . Using Theorem 1 and Corollary 1 we have

$$\begin{aligned} R_n^{\text{minimax}} &\geq R_{n, \mu}^{\text{Bayes}} \\ &\geq -\int_{\Theta} d\mu(\theta^*) \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-\frac{nh^2(\theta^*, \tilde{\theta})}{2}} \\ &= -\frac{1}{M} \sum_{i=1}^M \log \frac{1}{M} \sum_{j=1}^M e^{-\frac{nh^2(\theta_i, \theta_j)}{2}} \\ &\geq \log M - \log \left(1 + (M-1)e^{-\frac{n\epsilon^2}{2}} \right) \\ &\geq -\log \left(\frac{1}{M} + e^{-\frac{n\epsilon^2}{2}} \right). \end{aligned}$$

Since this holds for all ϵ , it follows that

$$R_n^{\minimax} \geq \sup_{\epsilon \geq 0} \left\{ -\log \left(\frac{1}{\mathcal{M}_\epsilon(\Theta, h)} + e^{-\frac{n\epsilon^2}{2}} \right) \right\}.$$

To complete the proof of part (1), simply note that $-\log(x + y) \geq -\log(2 \max(x, y)) = -\log 2 + \min\{-\log x, -\log y\}$. It follows that

$$R_n^{\minimax} \geq \sup_{\epsilon \geq 0} \min\{\log \mathcal{M}_\epsilon(\Theta, h), \frac{n\epsilon^2}{2}\} - \log 2.$$

Since $\mathcal{K}_{2\epsilon} = \log \mathcal{D}_{2\epsilon} \leq \log \mathcal{M}_\epsilon$, replacing ϵ with $\epsilon/2$, the second inequality follows.

We now turn to the upper bounds in part (2). Let $\Pi = \{\pi_1, \dots, \pi_M\}$ be any partition of Θ of diameter at most ϵ . For any prior measure μ on Θ , let $\mu_i = \mu(\pi_i)$. Then we use Theorem 1 and the upper bound given in Theorem 2 as follows.

$$\begin{aligned} R_n^{\minimax} &= \sup_{\mu} R_{n,\mu}^{\text{Bayes}} \\ &\leq \sup_{\mu} \left\{ -\int_{\Theta} d\mu(\theta^*) \log \int_{\Theta} d\mu(\tilde{\theta}) e^{-nD_{KL}(P_{\tilde{\theta}}||P_{\theta^*})} \right\} \\ &= \sup_{\mu} \left\{ -\sum_i \mu_i \int_{\pi_i} \frac{d\mu(\theta^*)}{\mu_i} \log \sum_j \mu_j \int_{\pi_j} \frac{d\mu(\tilde{\theta})}{\mu_j} e^{-nD_{KL}(P_{\tilde{\theta}}||P_{\theta^*})} \right\} \\ &\leq \sup_{\mu} \left\{ -\sum_i \mu_i \log \left(\mu_i e^{-b(\epsilon)n\epsilon^2} \right) \right\} \\ &= \sup_{\mu} \left\{ -\sum_i \mu_i \log \mu_i \right\} + b(\epsilon)n\epsilon^2 \\ &= \log M + b(\epsilon)n\epsilon^2. \end{aligned}$$

The second inequality follows by ignoring all but the i th term in the inner sum whenever the index on the outer sum is i , and noting that because the diameter of π_i is at most ϵ ,

$$D_{KL}(P_{\tilde{\theta}}||P_{\theta^*}) \leq b(\epsilon)h^2(\theta^*, \tilde{\theta}) \leq b(\epsilon)\epsilon^2$$

for all $\theta^*, \tilde{\theta} \in \pi_i$. The last equality follows from the fact that the entropy of a finite distribution is maximal for the uniform distribution. Since the particular partition of diameter ϵ can be chosen arbitrarily in the above chain of inequalities, it follows that $R_n^{\minimax} \leq \mathcal{K}_\epsilon(\Theta, h) + b(\epsilon)n\epsilon^2$ for any ϵ . This establishes the first inequality of part (2). The second inequality follows since $b(\epsilon) \leq b_{1/2}(\Theta)$ for all ϵ . The third inequality, but with $\sup_{\mu} R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}}$ in place of $R_{1,\rho_{1+\lambda}}^{\minimax}$, follows by an argument similar to that used for the first inequality, using Theorem 3. Since $\maximin \leq \minimax$ always, we have $\sup_{\mu} R_{1,\mu,\rho_{1+\lambda}}^{\text{Bayes}} \leq R_{1,\rho_{1+\lambda}}^{\minimax}$, and from this we obtain the result stated in the Theorem. \square

The method used in obtaining the upper bound in the above result is a familiar one (see e.g. [5, 34]). The method for obtaining the lower bound by choosing a discrete prior on a well-separated set of θ is also similar in many respects to standard lower bound methods, such as those that use Fano's inequality or Assouad's lemma (see e.g. [12, 10, 59]), but the

method is particularly clean in the present framework, giving a fairly good match to the upper bound.

In some cases \mathcal{K}_ϵ may not be a continuous function of ϵ , and even so it may not be obvious what kinds of asymptotic bounds on the risk $R_n^{minimax}$ are implied by Lemma 7. For such cases we make the following definitions.

Fix a totally bounded Θ and let $f_l(x)$ and $f_u(x)$ be any continuous, nondecreasing, unbounded functions on $(0, \infty)$ such that

$$\liminf_{\epsilon \rightarrow \infty} \frac{\mathcal{K}_\epsilon(\Theta, h)}{f_l(1/\epsilon)} \geq 1 \quad \text{and} \quad \limsup_{\epsilon \rightarrow \infty} \frac{\mathcal{K}_\epsilon(\Theta, h)}{f_u(1/\epsilon)} \leq 1. \quad (13)$$

For every positive real n let $\epsilon_l(n)$ be the unique solution to the equation $f_l(1/\epsilon) = n\epsilon^2$, and let $\epsilon_u(n)$ be the unique solution to the equation $f_u(1/\epsilon) = n\epsilon^2$. Let

$$F_l(n) = f_l\left(\frac{1}{\epsilon_l(n)}\right) = n\epsilon_l^2(n) \quad \text{and} \quad F_u(n) = f_u\left(\frac{1}{\epsilon_u(n)}\right) = n\epsilon_u^2(n). \quad (14)$$

Then we have the following lemma.

Lemma 8 *For every integer $n \geq 1$,*

1.

$$\liminf_{n \rightarrow \infty} \frac{R_n^{minimax}}{F_l(n/8)} \geq 1.$$

2. *If $\lim_{\epsilon \rightarrow 0} b(\epsilon) < \infty$ then for any function $h(n)$ such that $h(n) \rightarrow \infty$ as $n \rightarrow \infty$,*

$$\limsup_{n \rightarrow \infty} \frac{R_n^{minimax}}{F_u(nh(n))} \leq 1$$

and if there exists $\lambda > 0$ such that $R_{1, \rho_{1+\lambda}}^{minimax} < \infty$ then

$$\limsup_{n \rightarrow \infty} \frac{R_n^{minimax}}{F_u(nh(n) \log n)} \leq 1.$$

Proof: Using Lemma 7 and the definitions of f_l and F_l , we have

$$\begin{aligned} \liminf_{n \rightarrow \infty} \frac{R_n^{minimax}}{F_l(n/8)} &\geq \liminf_{n \rightarrow \infty} \frac{\min\left(\mathcal{K}_{\epsilon_l(n/8)}, \frac{n}{8} \epsilon_l^2\left(\frac{n}{8}\right)\right)}{F_l(n/8)} \\ &\geq \min\left(\liminf_{n \rightarrow \infty} \frac{\mathcal{K}_{\epsilon_l(n/8)}}{F_l(n/8)}, \liminf_{n \rightarrow \infty} \frac{\frac{n}{8} \epsilon_l^2\left(\frac{n}{8}\right)}{F_l(n/8)}\right) \\ &\geq \min\left(\liminf_{n \rightarrow \infty} \frac{f_l(1/\epsilon_l(n/8))}{F_l(n/8)}, 1\right) \\ &= 1. \end{aligned}$$

Now let $N = N(n) = nh(n)$. Let $\lim_{\epsilon \rightarrow 0} b(\epsilon) = b < \infty$. Then we also have

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{R_n^{minimax}}{F_u(nh(n))} &\leq \limsup_{n \rightarrow \infty} \frac{\mathcal{K}_{\epsilon_u(N)} + bn\epsilon_u^2(N)}{F_u(N)} \\ &\leq \limsup_{n \rightarrow \infty} \left(\frac{f_u(1/\epsilon_u(N))}{F_u(N)} + \frac{b}{h(n)} \right) \\ &= 1 + \limsup_{n \rightarrow \infty} \frac{b}{h(n)} \\ &= 1. \end{aligned}$$

The last inequality follows similarly, using the last inequality of Lemma 7. \square

Essentially, when $F_l(n)$ and $F_u(n \log n)$ are close asymptotically, as can often be arranged, this lemma shows that asymptotic rates for $R_n^{minimax}$ can be obtained by “solving” the equation $\mathcal{K}_\epsilon(\Theta, h) = n\epsilon^2$. This general approach was developed by Le Cam and Birgé [42, 11, 12]. We illustrate it by applying the above lemma to all of the standard cases for the asymptotic growth rate of the metric entropy $\mathcal{K}_\epsilon(\Theta, h)$.

Theorem 4 *Assume there exists $\lambda > 0$ such that $R_{1, \rho_{1+\lambda}}^{minimax} < \infty$.⁵*

1. *If Θ is finite then*

$$R_n^{minimax} \rightarrow \log |\Theta| \text{ as } n \rightarrow \infty.$$

2. *If $\dim(\Theta, h) = 0$ then*

$$R_n^{minimax} \in o(\log n).$$

3. *If $\dim(\Theta, h) = D$ where $0 < D < \infty$ then*

$$R_n^{minimax} \sim \frac{D}{2} \log n.$$

4. *If $\text{df}(\Theta, h) = \beta$ where $1 < \beta < \infty$ then*

$$\log R_n^{minimax} \sim \beta \log \log n.$$

5. *If $\text{mo}(\Theta, h) = \alpha$ where $0 < \alpha < \infty$ then*

$$\log R_n^{minimax} \sim \frac{\alpha}{2 + \alpha} \log n.$$

6. *If $\text{mo}(\Theta, h) = \infty$ or (Θ, h) is not totally bounded, then*

$$\text{if } R_1^{minimax} < \infty \text{ then } \log R_n^{minimax} \sim \log n \text{ else } R_n^{minimax} = \infty \text{ for all } n.$$

⁵Actually, only the upper bounds in parts (2)-(5) require the assumption that there exists $\lambda > 0$ such that $R_{1, \rho_{1+\lambda}}^{minimax} < \infty$.

Proof: As mentioned after Corollary 2, part (1) follows from that Corollary and Theorem 1. Each of the results (2)-(5) follows easily from Lemma 8 by plugging in the appropriate rates for f_l and f_u , and solving for F_l and F_u . We illustrate this for parts (3) and (5); the other parts are similar. For part (3), since $\dim(\Theta, h) = D$ where $0 < D < \infty$, we may choose

$$f_l(x) = f_u(x) = D \log x.$$

Solving $D \log \frac{1}{\epsilon} = n\epsilon^2$, we find that

$$\epsilon_l(n) = \epsilon_u(n) \sim \sqrt{\frac{D}{2n} \log n},$$

and hence by (14)

$$F_l(n) = F_u(n) \sim \frac{D}{2} \log n.$$

From the lower bound of Lemma 8, it follows that

$$\liminf_{n \rightarrow \infty} \frac{R_n^{minimax}}{\frac{D}{2} \log \frac{n}{8}} \geq 1.$$

Let $h(n) = \log n$. From the second upper bound of Lemma 8, it follows that

$$\limsup_{n \rightarrow \infty} \frac{R_n^{minimax}}{\frac{D}{2} \log(n \log^2 n)} \leq 1.$$

The result in part (3) follows.

In part (5), since $\mathbf{mo}(\Theta, h) = \alpha$ where $0 < \alpha < \infty$, for any $0 < \delta < \alpha$ we can choose $f_l(x) = x^{\alpha-\delta}$ and $f_u(x) = x^{\alpha+\delta}$. Solving $\epsilon^{\alpha \pm \delta} = n\epsilon^2$, we find that

$$F_l(n) = n^{\frac{\alpha-\delta}{2+\alpha-\delta}} \quad \text{and} \quad F_u(n) = n^{\frac{\alpha+\delta}{2+\alpha+\delta}}$$

From the lower bound of Lemma 8, it follows that for all $0 < \delta < \alpha$,

$$\liminf_{n \rightarrow \infty} \frac{R_n^{minimax}}{(n/8)^{\frac{\alpha-\delta}{2+\alpha-\delta}}} \geq 1.$$

Hence

$$\liminf_{n \rightarrow \infty} \frac{\log R_n^{minimax}}{\frac{\alpha}{2+\alpha} \log n} \geq 1.$$

Now let $h(n) = \log n$. Then from the second upper bound of Lemma 8, it follows that for all $0 < \delta < \alpha$,

$$\limsup_{n \rightarrow \infty} \frac{R_n^{minimax}}{(n \log^2 n)^{\frac{\alpha+\delta}{2+\alpha+\delta}}} \leq 1.$$

Hence

$$\limsup_{n \rightarrow \infty} \frac{\log R_n^{minimax}}{\frac{\alpha}{2+\alpha} \log n} \leq 1.$$

This establishes part (5).

To verify part (6), first note that the minimax risk $R_n^{minimax}$ is nondecreasing in n . Furthermore, if $R_n^{minimax}$ is finite, then it can grow at most linearly, as is seen in the following series of inequalities.

$$\begin{aligned}
R_n^{minimax} &= \inf_{\text{dist. } R \text{ on } Y^n} \sup_{\theta \in \Theta} D_{KL}(P_\theta^n || R) \\
&\leq \inf_{\text{dist. } Q \text{ on } Y} \sup_{\theta \in \Theta} D_{KL}(P_\theta^n || Q^n) \\
&= n \inf_{\text{dist. } Q \text{ on } Y} \sup_{\theta \in \Theta} D_{KL}(P_\theta || Q) \\
&= n R_1^{minimax}
\end{aligned}$$

Hence for any Θ , either $R_n^{minimax} = \infty$ for all n or $R_n^{minimax}$ is finite and bounded by $n R_1^{minimax}$ for all n .

If (Θ, h) is not totally bounded then $\mathcal{M}_{\epsilon_0}(\Theta, h)$ is infinite for some $\epsilon_0 > 0$. In this case the first lower bound from part (1) of Lemma 7 shows that

$$R_n^{minimax} \geq \frac{n \epsilon_0^2}{2}$$

for all $n \geq 1$. Hence $R_n^{minimax} \asymp n$ in this case. If (Θ, h) is totally bounded but $\mathbf{mo}(\Theta, h) = \infty$ then $R_1^{minimax} < \infty$ and by the same reasoning as in the proof of part (5), for all $\alpha > 0$,

$$\liminf_{n \rightarrow \infty} \frac{R_n^{minimax}}{(n/8)^{\frac{\alpha}{2+\alpha}}} \geq 1.$$

This, combined with the fact that $R_n^{minimax} \leq n R_1^{minimax}$, implies that $\log R_n^{minimax} \sim \log n$. \square

The above theorem does not give very precise bounds in the infinite dimensional case. Indeed, not much more can be said using the fairly crude notions of metric order and functional dimension to measure the massiveness of infinite dimensional Θ . To remedy this, below is a more refined result.

Theorem 5 *Assume there exists $\lambda > 0$ such that $R_{1, \rho_{1+\lambda}}^{minimax} < \infty$. Let $l(x)$ be a continuous, nondecreasing function defined on the positive reals such that for all $\gamma \geq 0$ and $C > 0$*

1.

$$\lim_{x \rightarrow \infty} \frac{l(Cx(l(x))^\gamma)}{l(x)} = 1$$

and

2.

$$\lim_{x \rightarrow \infty} \frac{l(Cx(\log(x))^\gamma)}{l(x)} = 1.$$

Then

1.

$$\text{If } \mathcal{K}_\epsilon(\Theta, h) \sim l\left(\frac{1}{\epsilon}\right) \text{ then } R_n^{\text{minimax}} \sim l(\sqrt{n}).$$

2. If for some $\alpha > 0$,

$$\mathcal{K}_\epsilon(\Theta, h) \asymp \left(\frac{1}{\epsilon}\right)^\alpha l\left(\frac{1}{\epsilon}\right)$$

then

(a)

$$\text{If } \lim_{\epsilon \rightarrow 0} b(\epsilon) < \infty \text{ then } R_n^{\text{minimax}} \asymp n^{\alpha/(\alpha+2)} \left[l(n^{1/(\alpha+2)})\right]^{2/(\alpha+2)}$$

else

(b)

$$\liminf_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{n^{\alpha/(\alpha+2)} \left[l(n^{1/(\alpha+2)})\right]^{2/(\alpha+2)}} > 0$$

and

$$\limsup_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{n^{\alpha/(\alpha+2)} \left[l(n^{1/(\alpha+2)})\right]^{2/(\alpha+2)} (\log n)^{\alpha/(\alpha+2)}} < \infty.$$

Proof: Consider part (2) first. Since $\mathcal{K}_\epsilon(\Theta, h) \asymp \left(\frac{1}{\epsilon}\right)^\alpha l\left(\frac{1}{\epsilon}\right)$ we may choose

$$f_l(x) = ax^\alpha l(x) \text{ and } f_u(x) = bx^\alpha l(x)$$

for suitable constants $0 < a \leq b$. Solving $f_l(x) = n/x^2$, we find that

$$x \sim \left(\frac{N}{l(N^{1/(\alpha+2)})}\right)^{1/(\alpha+2)},$$

where $N = n/a$. Here we use property (1) of $l(x)$. Hence

$$\epsilon_l(n) \sim \left(\frac{l(N^{1/(\alpha+2)})}{N}\right)^{1/(\alpha+2)},$$

and thus by (14), and again using property (1) of $l(x)$,

$$F_l(n) \asymp n^{\alpha/(\alpha+2)} \left[l(n^{1/(\alpha+2)})\right]^{2/(\alpha+2)}.$$

By similar reasoning

$$F_u(n) \asymp n^{\alpha/(\alpha+2)} \left[l(n^{1/(\alpha+2)})\right]^{2/(\alpha+2)}.$$

From the lower bound of Lemma 8, and property (1), it follows that

$$\liminf_{n \rightarrow \infty} \frac{R_n^{\text{minimax}}}{n^{\alpha/(\alpha+2)} \left[l(n^{1/(\alpha+2)})\right]^{2/(\alpha+2)}} > 0$$

From the second upper bound of Lemma 8, it follows that for any unbounded, increasing function $g(n)$,

$$\limsup_{n \rightarrow \infty} \frac{R_n^{minimax}}{n^{\alpha/(\alpha+2)} [l((ng(n) \log n)^{1/(\alpha+2)})]^{2/(\alpha+2)} (g(n) \log n)^{\alpha/(\alpha+2)}} < \infty.$$

Part (2b) follows easily from this, using property (2) of the function $l(x)$. For part (2a), note that from the first upper bound of Lemma 8, if $\lim_{\epsilon \rightarrow 0} b(\epsilon) < \infty$ then the $\log n$ factors can be removed from the lim sup above, yielding the desired result. Part (1) follows by a similar argument, essentially setting $\alpha = 0$, and $a = b = 1$, so that most terms in the denominators of the expressions above go away, and tracking the lim inf and lim sup more precisely. \square

Note that in finite dimensional cases, we have $\mathcal{K}_\epsilon \sim D \log \frac{1}{\epsilon} = l(\frac{1}{\epsilon})$, and part (1) of the above theorem gives $R_n^{minimax} \sim l(\sqrt{n}) = \frac{D}{2} \log n$, as obtained in the previous theorem. Part (1) generalizes this to infinite dimensional cases of finite functional dimension, in which, e.g., $\mathcal{K}_\epsilon \sim C \left(\log \frac{1}{\epsilon}\right)^\beta$ for $\beta > 1$. Part (2) does the same for cases in which Θ has finite metric order.

The above theorem is not applicable in all cases. In particular, it can be shown that the condition that $R_{1, \rho_{1+\lambda}}^{minimax} < \infty$ in the above result and the preceding results of this section cannot be removed. For example, this condition is violated by the Θ defined in Example 1. In this case (Θ, h) is totally bounded and $R_n^{minimax} \sim n$, yet $\mathcal{K}_\epsilon \sim (1/\epsilon)^2$, which would yield via Theorem 5 an estimated rate of \sqrt{n} for $R_n^{minimax}$. This is off by a factor of \sqrt{n} . Of course the lower bounds in Theorem 5 and the preceding results are valid in this and any other case without any special assumptions, but in this case, we see that they are not tight.

8 Bounds on instantaneous minimax risk for various loss functions

Here we show how the results of the previous sections can be used to give upper and lower bounds on the instantaneous minimax risk of estimating a probability distribution for various loss functions, as defined in section 3. One way to do this is to use Fano's inequality, as described in [10]. Here we look at a simple alternate approach.

Recall that the instantaneous minimax risk at time t is denoted

$$r_{t,L}^{minimax} = \inf_{\hat{P}} \sup_{\theta^* \in \Theta} \int_{Y^{t-1}} dP_{\theta^*}^{t-1} L(P_{\theta^*}, \hat{P}_t),$$

for loss function L . The instantaneous Bayes risk under prior μ is defined similarly by taking expectation over Θ^* instead of supremum, and denoted $r_{t,\mu,L}^{Bayes}$. Below we will consider other loss functions, but for now we assume that $L(P, Q) = D_{KL}(P||Q)$ and omit the subscript L .

While it is easily verified that $R_{n,\mu}^{Bayes} = \sum_{t=1}^n r_{t,\mu}^{Bayes}$, the exact relationship between the instantaneous and cumulative minimax risks is less clear. However, Barron et al. [5, 18, 10] have shown the following.

Lemma 9 [5]

$$\sum_{t=1}^n r_t^{minimax} \geq R_n^{minimax} \geq n r_n^{minimax}$$

Proof: For the first inequality, simply note that

$$\begin{aligned}
\sum_{t=1}^n r_t^{minimax} &= \sum_{t=1}^n \inf_{\hat{P}} \sup_{\theta^* \in \Theta} \int_{Y^{t-1}} dP_{\theta^*}^{t-1} D_{KL}(P_{\theta^*} || \hat{P}_t) \\
&= \inf_{\hat{P}} \sum_{t=1}^n \sup_{\theta^* \in \Theta} \int_{Y^{t-1}} dP_{\theta^*}^{t-1} D_{KL}(P_{\theta^*} || \hat{P}_t) \\
&\geq \inf_{\hat{P}} \sup_{\theta^* \in \Theta} \sum_{t=1}^n \int_{Y^{t-1}} dP_{\theta^*}^{t-1} D_{KL}(P_{\theta^*} || \hat{P}_t) \\
&= R_n^{minimax}
\end{aligned}$$

For the second inequality, let μ be any prior on Θ . Let $M_{n,\mu} = \int_{\Theta} P_{\theta}^n d\mu(\theta)$ be the Bayes mixture for μ , and $P_{t,\mu}^{Bayes}(y_t|y^{t-1})$ be the posterior predictive distribution for each $1 \leq t \leq n$. Fix n . For each y^{n-1} , let the strategy Q_{μ} be defined by the predictive distributions

$$Q_{n,\mu} = \frac{1}{n} \sum_{t=1}^n P_{t,\mu}^{Bayes}.$$

Then for all $\theta^* \in \Theta$, the instantaneous risk of the strategy Q_{μ} at time n is

$$\begin{aligned}
r_{n,Q_{\mu}}(\theta^*) &= \int_{Y^{n-1}} dP_{\theta^*}^{n-1} D_{KL}(P_{\theta^*} || Q_{n,\mu}) \\
&\leq \frac{1}{n} \sum_{t=1}^n \int_{Y^{n-1}} dP_{\theta^*}^{n-1} D_{KL}(P_{\theta^*} || P_{t,\mu}^{Bayes}) \\
&= \frac{1}{n} D_{KL}(P_{\theta^*}^n || M_{n,\mu}) \\
&= \frac{R_{n,P_{\mu}^{Bayes}}(\theta^*)}{n},
\end{aligned}$$

where P_{μ}^{Bayes} is the Bayes strategy for the cumulative risk under prior μ . Here the inequality follows from Jensen's inequality and the next equality from the chain rule for relative entropy (see e.g. [21], p. 23). It follows that the instantaneous minimax risk at time n is

$$\begin{aligned}
r_n^{minimax} &= \inf_{\hat{P}} \sup_{\theta^* \in \Theta} r_{n,\hat{P}}(\theta^*) \\
&\leq \inf_{\mu} \sup_{\theta^* \in \Theta} r_{n,Q_{\mu}}(\theta^*) \\
&\leq \frac{1}{n} \inf_{\mu} \sup_{\theta^* \in \Theta} R_{n,P_{\mu}^{Bayes}}(\theta^*) \\
&= \frac{1}{n} \inf_{\hat{P}} \sup_{\theta^* \in \Theta} R_{n,\hat{P}}(\theta^*) \\
&= \frac{R_n^{minimax}}{n}.
\end{aligned}$$

The penultimate equality follows the second part of Theorem 1. \square

Using Lemma 8, the above lemma may be further refined to give the following relationships between the cumulative and instantaneous minimax risks under relative entropy loss.

Lemma 10 Fix a totally bounded Θ and let $f_l(x)$ and $f_u(x)$ be any continuous, nondecreasing, unbounded functions on $(0, \infty)$ such that

$$\liminf_{\epsilon \rightarrow \infty} \frac{\mathcal{K}_\epsilon(\Theta, h)}{f_l(1/\epsilon)} \geq 1 \quad \text{and} \quad \limsup_{\epsilon \rightarrow \infty} \frac{\mathcal{K}_\epsilon(\Theta, h)}{f_u(1/\epsilon)} \leq 1. \quad (15)$$

For every positive real n let $\epsilon_l(n)$ be the unique solution to the equation $f_l(1/\epsilon) = n\epsilon^2$, and let $\epsilon_u(n)$ be the unique solution to the equation $f_u(1/\epsilon) = n\epsilon^2$. Let

$$F_l(n) = f_l\left(\frac{1}{\epsilon_l(n)}\right) = n\epsilon_l^2(n) \quad \text{and} \quad F_u(n) = f_u\left(\frac{1}{\epsilon_u(n)}\right) = n\epsilon_u^2(n). \quad (16)$$

Then for every integer $n \geq 1$:

1. If $r_t^{minimax} < \infty$ for all t , F_l is differentiable and its derivative F_l' is nonincreasing, then

$$\limsup_{n \rightarrow \infty} \frac{8r_n^{minimax}}{F_l'(n/8)} \geq 1.$$

2. If $\lim_{\epsilon \rightarrow 0} b(\epsilon) < \infty$ then for any function $h(n)$ such that $h(n) \rightarrow \infty$ as $n \rightarrow \infty$,

$$\limsup_{n \rightarrow \infty} \frac{nr_n^{minimax}}{F_u(nh(n))} \leq 1$$

and if there exists $\lambda > 0$ such that $R_{1, \rho_{1+\lambda}}^{minimax} < \infty$ then

$$\limsup_{n \rightarrow \infty} \frac{nr_n^{minimax}}{F_u(nh(n) \log n)} \leq 1.$$

Proof: The upper bounds follow directly from part (2) of Lemma 8, using the above result that $R_n^{minimax} \geq nr_n^{minimax}$. For the lower bound, let $G(n) = F_l(n/8)$ and $g(n) = G'(n) = \frac{1}{8}F_l'(n/8)$. From part (1) of Lemma 8, and the above result that $\sum_{t=1}^n r_t^{minimax} \geq R_n^{minimax}$ we have

$$\liminf_{n \rightarrow \infty} \frac{\sum_{t=1}^n r_t^{minimax}}{G(n)} \geq 1.$$

By the fundamental theorem of calculus $G(n) = \int_0^n g(t)dt + G(0)$. Since g is nonincreasing, $\int_0^n g(t)dt \geq \sum_{t=1}^n g(t)$. Since f_l is unbounded, G is unbounded, and thus so is $\sum_{t=1}^n g(t)$. It follows that

$$\liminf_{n \rightarrow \infty} \frac{\sum_{t=1}^n r_t^{minimax}}{\sum_{t=1}^n g(t)} \geq 1.$$

Again, since $\sum_{t=1}^n g(t)$ is unbounded, and since $r_t^{minimax} < \infty$ for all t , this implies that

$$\limsup_{n \rightarrow \infty} \frac{r_n^{minimax}}{g(n)} \geq 1.$$

This gives the result. \square

In order to use the above lemma for other loss functions, such as squared Hellinger and L_1 loss functions, we need only bound these loss functions in terms of the relative entropy loss. The following bounds are well known for any distributions P and Q on Y (see e.g. [42])

$$D_{HL}^2(P, Q) \leq \|P - Q\| \leq 2D_{HL}(P, Q) \quad (17)$$

In addition, using Lemma 4, it follows that

$$D_{HL}^2(P, Q) \leq D_{KL}(P||Q) \leq \sup_{y \in Y} b_{1/2} \left(\frac{dP(y)}{dQ(y)} \right) D_{HL}^2(P, Q) \quad (18)$$

Finally, using Lemma 5 with $\epsilon = \epsilon_n = \frac{1}{n^{2/\lambda} \log n}$, it can easily be shown that for any distribution U such that $C = \int (dP)^{1+\lambda} (dU)^{-\lambda} < \infty$, there exists $N \geq 1$ such that for all $n \geq N$,

$$D_{KL}(P||((1 - \epsilon_n)Q + \epsilon_n U)) \leq \frac{8 \log n}{\lambda} D_{HL}^2(P, Q) + \frac{2C}{n \log^{\lambda/2} n} \quad (19)$$

Let $r_{n, D_{HL}^2}^{minimax}$ denote the instantaneous minimax risk for the squared Hellinger loss, and $r_{n, \|\cdot\|}^{minimax}$ denote the instantaneous minimax risk for the L_1 loss. The above inequalities imply the following relationships between these risks and the instantaneous minimax risk under relative entropy loss, $r_n^{minimax}$.

Lemma 11 1.

$$\frac{r_n^{minimax}}{b_{1/2}(\Theta)} \leq r_{n, D_{HL}^2}^{minimax} \leq r_n^{minimax}$$

2.

$$r_{n, D_{HL}^2}^{minimax} \leq r_{n, \|\cdot\|}^{minimax} \leq 2\sqrt{r_{n, D_{HL}^2}^{minimax}}$$

3. If there exists $0 \leq \lambda \leq 1$ such that $R_{1, \rho_{1+\lambda}}^{minimax} < \infty$ then there exists an $N \geq 1$ such that for all $n \geq N$,

$$r_n^{minimax} \leq \frac{8 \log n}{\lambda} r_{n, D_{HL}^2}^{minimax} + \frac{2R_{1, \rho_{1+\lambda}}^{minimax}}{n \log^{\lambda/2} n}$$

Proof: Most of these results follow directly from the corresponding inequalities above. Only two of them require comment. For the second inequality in part (2), note that

$$\begin{aligned} r_{n, \|\cdot\|}^{minimax} &= \inf_{\hat{P}} \sup_{\theta^* \in \Theta} \int_{Y^{n-1}} dP_{\theta^*}^{n-1} \|P_{\theta^*} - \hat{P}_n\| \\ &\leq 2 \inf_{\hat{P}} \sup_{\theta^* \in \Theta} \int_{Y^{n-1}} dP_{\theta^*}^{n-1} D_{HL}(P_{\theta^*}, \hat{P}_n) \\ &\leq 2 \inf_{\hat{P}} \sup_{\theta^* \in \Theta} \sqrt{\int_{Y^{n-1}} dP_{\theta^*}^{n-1} D_{HL}^2(P_{\theta^*}, \hat{P}_n)} \\ &= 2 \sqrt{\inf_{\hat{P}} \sup_{\theta^* \in \Theta} \int_{Y^{n-1}} dP_{\theta^*}^{n-1} D_{HL}^2(P_{\theta^*}, \hat{P}_n)} \\ &= 2 \sqrt{r_{n, D_{HL}^2}^{minimax}}. \end{aligned}$$

The first inequality comes from (17) and the second from Jensen's inequality.

To see the inequality in part (3), for any $\gamma > 0$, let U_γ be a distribution such that

$$C_\gamma = \sup_{\theta \in \Theta} \int (dP_\theta)^{1+\lambda} (dU_\gamma)^{-\lambda} \leq R_{1,\rho_{1+\lambda}}^{minimax} + \gamma.$$

Then $\inf_{\gamma > 0} C_\gamma = R_{1,\rho_{1+\lambda}}^{minimax}$. Now note that for all $\gamma > 0$,

$$\begin{aligned} r_n^{minimax} &= \inf_{\hat{P}} \sup_{\theta^* \in \Theta} \int_{Y^{n-1}} dP_{\theta^*}^{n-1} D_{KL}(P_{\theta^*} || \hat{P}_n) \\ &\leq \inf_{\hat{P}} \sup_{\theta^* \in \Theta} \int_{Y^{n-1}} dP_{\theta^*}^{n-1} D_{KL}(P_{\theta^*} || (1 - \epsilon_n)\hat{P}_n + \epsilon_n U_\gamma) \\ &\leq \frac{8 \log n}{\lambda} \inf_{\hat{P}} \sup_{\theta^* \in \Theta} \int_{Y^{n-1}} dP_{\theta^*}^{n-1} D_{HL}^2(P_{\theta^*}, \hat{P}_n) + \frac{2C_\gamma}{n \log^{\lambda/2} n} \\ &= \frac{8 \log n}{\lambda} r_{n,D_{HL}^2}^{minimax} + \frac{2C_\gamma}{n \log^{\lambda/2} n}. \end{aligned}$$

The result follows. \square

We illustrate this lemma first by giving a simple proof that when Θ is finite, the instantaneous minimax risk decreases exponentially in n for all of the above loss functions.

Theorem 6 *If Θ is finite then there exist $a > 0$ and $0 < b < 1$ such that for all $n \geq 1$,*

$$r_n^{minimax}, r_{n,D_{HL}^2}^{minimax} \leq ab^n \quad \text{and} \quad r_{n,||\cdot||}^{minimax} \leq ab^{n/2}.$$

Proof: Let μ be the uniform distribution on Θ . By Corollary 3, there exist $A > 0$ and $0 < b < 1$ such that for all n

$$R_{n,\mu}^{Bayes} \geq \log |\Theta| - Ab^n.$$

Since $R_{n,\mu}^{Bayes} = \sum_{t=1}^n r_{t,\mu}^{Bayes}$, it follows that

$$r_{n,\mu}^{Bayes} \leq \sum_{t=n}^{\infty} r_{t,\mu}^{Bayes} = \log |\Theta| - R_{n-1,\mu}^{Bayes} \leq Ab^{n-1}.$$

Thus if $\hat{P} = P_\mu^{Bayes}$, the Bayes strategy for the uniform prior μ , then

$$\sum_{\theta^* \in \Theta} \frac{r_{n,\hat{P}}(\theta^*)}{|\Theta|} \leq Ab^{n-1}.$$

It follows that for all $\theta^* \in \Theta$

$$r_{n,\hat{P}}(\theta^*) \leq |\Theta| Ab^{n-1}.$$

Let $a = |\Theta|Ab$. Then

$$r_n^{minimax} \leq r_{n,\hat{P}} \leq ab^n.$$

The remaining inequalities follow from the bounds in parts (1) and (2) in Lemma 11, respectively. For the second bound we can replace the previous a with $2\sqrt{a}$. \square

In order to see how these results can be used to derive bounds on the instantaneous minimax risk for more general Θ , let us make the following definition. For any loss L , let us define the *best exponent for the instantaneous minimax risk* by

$$e_L = \sup\{x : \limsup_{n \rightarrow \infty} \frac{r_{n,L}^{minimax}}{n^{-x}} \leq 1\}.$$

Theorem 7 *Assume there exists $\lambda > 0$ such that $R_{1,\rho_{1+\lambda}}^{minimax} < \infty$. Then the bounds on e_L given in the following table are valid.*⁶

size of Θ	loss function		
	D_{KL}	D_{HL}^2	$\ \cdot\ $
Θ is finite	$e_{D_{KL}} = \infty$	$e_{D_{HL}^2} = \infty$	$e_{\ \cdot\ } = \infty$
$\mathbf{dim}(\Theta, h) = 0$	$e_{D_{KL}} \geq 1$	$e_{D_{HL}^2} \geq 1$	$e_{\ \cdot\ } \geq 1/2$
$\mathbf{dim}(\Theta, h) = D$ where $0 < D < \infty$	$e_{D_{KL}} = 1$	$e_{D_{HL}^2} = 1$	$1/2 \leq e_{\ \cdot\ } \leq 1$
$\mathbf{df}(\Theta, h) = \beta$ where $1 < \beta < \infty$	$e_{D_{KL}} = 1$	$e_{D_{HL}^2} = 1$	$1/2 \leq e_{\ \cdot\ } \leq 1$
$\mathbf{mo}(\Theta, h) = \alpha$ where $0 < \alpha < \infty$	$e_{D_{KL}} = \frac{2}{2+\alpha}$	$e_{D_{HL}^2} = \frac{2}{2+\alpha}$	$\frac{1}{2+\alpha} \leq e_{\ \cdot\ } \leq \frac{2}{2+\alpha}$
$\mathbf{mo}(\Theta, h) = \infty$	$e_{D_{KL}} = 0$	$e_{D_{HL}^2} = 0$	$e_{\ \cdot\ } = 0$
(Θ, h) not totally bounded	$e_{D_{KL}} = 0$	$e_{D_{HL}^2} = 0$	$e_{\ \cdot\ } = 0$

Proof: The results for finite Θ follow directly from the previous theorem. The remaining results follow from Lemma 10 and Lemma 11, and in each case the proof is analogous to that of the corresponding result in Theorem 4. We give the derivation in two of the cases; the remaining derivations are similar.

As in the proof of Theorem 4, if $\mathbf{dim}(\Theta, h) = D$ where $0 < D < \infty$, we may choose $f_l(x) = f_u(x) = D \log x$, and hence

$$F_l(n) = F_u(n) \sim \frac{D}{2} \log n.$$

From the lower bound of Lemma 10, it follows that

$$\limsup_{n \rightarrow \infty} \frac{16nr_n^{minimax}}{D} \geq 1.$$

Let $h(n) = \log n$. From the second upper bound of Lemma 10, it follows that

$$\limsup_{n \rightarrow \infty} \frac{2nr_n^{minimax}}{D \log(n \log^2 n)} \leq 1.$$

Hence $e_{D_{KL}} = 1$.

Since $r_{n,D_{HL}^2}^{minimax} \leq r_n^{minimax}$, from the upper bound above we also have

$$\limsup_{n \rightarrow \infty} \frac{2nr_{n,D_{HL}^2}^{minimax}}{D \log(n \log^2 n)} \leq 1.$$

⁶Actually, the upper bounds in the first column only require the weaker assumption that $r_t^{minimax} < \infty$ for all t , and in the case that Θ is finite none of the results require any additional assumptions.

Combining the lower bound above with part (3) of Lemma 11, we see that

$$\limsup_{n \rightarrow \infty} \frac{(128n \log n) r_{n, D_{HL}^2}^{minimax}}{\lambda D} + \limsup_{n \rightarrow \infty} \frac{32 R_{1, \rho_{1+\lambda}}^{minimax}}{D \log^{\lambda/2} n} = \limsup_{n \rightarrow \infty} \frac{(128n \log n) r_{n, D_{HL}^2}^{minimax}}{\lambda D} \geq 1.$$

Hence $e_{D_{HL}^2} = 1$.

Since $r_{n, D_{HL}^2}^{minimax} \leq r_{n, \|\cdot\|}^{minimax}$, it follows that

$$\limsup_{n \rightarrow \infty} \frac{(128n \log n) r_{n, \|\cdot\|}^{minimax}}{\lambda D} \geq 1.$$

Since $r_{n, \|\cdot\|}^{minimax} \leq 2\sqrt{r_{n, D_{HL}^2}^{minimax}}$, from the above upper bound for $r_{n, D_{HL}^2}^{minimax}$ we get

$$\limsup_{n \rightarrow \infty} \frac{1}{2} \left(\sqrt{\frac{2n}{D \log(n \log^2 n)}} \right) r_{n, \|\cdot\|}^{minimax} \leq 1.$$

It follows that $1/2 \leq e_{\|\cdot\|} \leq 1$.

Skipping to the last line of the table, when $\mathbf{mo}(\Theta, h) = \alpha$ where $0 < \alpha < \infty$, for any $0 < \delta < \alpha$ we can choose $f_l(x) = x^{\alpha-\delta}$ and $f_u(x) = x^{\alpha+\delta}$, yielding

$$F_l(n) = n^{\frac{\alpha-\delta}{2+\alpha-\delta}} \quad \text{and} \quad F_u(n) = n^{\frac{\alpha+\delta}{2+\alpha+\delta}}$$

From the lower bound of Lemma 10, it follows that for all $0 < \delta < \alpha$,

$$\limsup_{n \rightarrow \infty} \frac{8r_n^{minimax}}{\frac{\alpha-\delta}{2+\alpha+\delta} (n/8)^{-\frac{2}{2+\alpha-\delta}}} \geq 1.$$

Now let $h(n) = \log n$. Then from the second upper bound of Lemma 10, it follows that for all $0 < \delta < \alpha$,

$$\limsup_{n \rightarrow \infty} \frac{nr_n^{minimax}}{(n \log^2 n)^{\frac{\alpha+\delta}{2+\alpha+\delta}}} \leq 1.$$

Hence $e_{D_{KL}} = \frac{2}{2+\alpha}$. The bounds for $e_{D_{HL}^2}$ and $e_{\|\cdot\|}$ in this case are then derived from the above bounds in the same manner they were previously for the finite dimensional case. \square

Again, it is easy to see the assumption that $R_{1, \rho_{1+\lambda}}^{minimax} < \infty$ cannot be removed from this theorem. In particular, for the Θ of Example 1, it can be shown that $\mathbf{mo}(\Theta, h) = 2$, yet $e_{D_{HL}^2}, e_{\|\cdot\|} \geq 1$ (too high for Theorem 7 to apply) and $e_{D_{KL}} = 0$ (too low).

9 Example application of the minimax bounds: Non-parametric density estimation

Here we give a brief, fairly classical example just to illustrate how one may apply the results given above. Let us assume that the statistician observes a set of n observations y_1, \dots, y_n which are drawn independently from a density $dP_\theta(y)$, $\theta \in \Theta$ on the interval $[0, 1]$. As in

[59], let Θ be the Lipschitz class $F_{p,\alpha}(C, L)$ of densities satisfying $\sup_{y \in [0,1]} |dP_\theta(y)| \leq C$ and having derivatives $dP_\theta^{(k)}(y)$ of order $k \leq p$ with the Lipschitz condition on the p -th derivative $|dP_\theta^{(p)}(y) - dP_\theta^{(p)}(y')| \leq L|y - y'|^\alpha$ for $y, y' \in [0, 1]$. Since the functions in $F_{p,\alpha}(C_1, L)$ are uniformly bounded, they have an integrable envelope function, and hence $R_{1,\rho_{1+\lambda}}^{minimax} < \infty$ for all $\lambda > 0$. Furthermore, since the functions in $F_{p,\alpha}(C_1, L)$ are uniformly bounded, all L_q distances ($q \geq 1$) are equivalent. As shown by Barron and Yang [10], a further restriction to uniformly lower bounded densities also insures that the condition $\lim_{\epsilon \rightarrow 0} b(\epsilon) < \infty$ holds, and makes the Hellinger distance equivalent to the L_q distances, without changing the metric entropy asymptotically. By a result of Clements [20], the metric entropy of Θ under L_1 distance is given by

$$\mathcal{K}_\epsilon(\Theta, L_1) \asymp \epsilon^{-\frac{1}{p+\alpha}}.$$

Hence

$$\mathcal{K}_\epsilon(\Theta, h) \asymp \epsilon^{-\frac{1}{p+\alpha}}.$$

Thus from Theorem 5 we get $R_n^{minimax} \asymp n^{\frac{1}{2(p+\alpha)+1}}$, and from Theorem 7 we get that the best exponent for the instantaneous minimax risk is $e_{DKL} = \frac{2(p+\alpha)}{2(p+\alpha)+1}$ when the loss is the KL-divergence, and the same when the loss is the squared Hellinger distance, while it is within a factor between 1/2 and 1 of this for the L_1 loss. Sharper results are known (see e.g. [12, 32]).

Since the metric entropies are known for many interesting classes of functions, many more examples of this type are possible. Many such examples are given by Birgé [11, 12] and Barron and Yang [10].

10 Discussion

We have shown that under relatively weak assumptions, (in particular, whenever there exists a distribution U and a $\lambda > 0$ such that the $(1 + \lambda)$ -affinity between P_θ and U is uniformly bounded for all $\theta \in \Theta$) one can obtain explicit bounds on the mutual information $I(\Theta^*; Y^n)$ between the true parameter and the observations in terms of a Laplace transform of the Hellinger distance in Θ , and from these one can obtain bounds on the cumulative minimax risk in estimating a distribution in Θ under relative entropy loss in terms of the metric entropy of Θ with respect to Hellinger distance. In fact, in each case only the upper bounds depend on the assumptions; the lower bounds hold for any Θ . We also show by example that some assumptions are needed to get the type of general characterizations of the mutual information and minimax risk in terms of the Hellinger distance that we obtain. It remains open to get a useful characterization of these quantities for the cases where our assumptions do not hold.

We also show how general bounds on instantaneous risk in estimating a distribution for various loss functions can be derived in a very simple manner from the bounds on cumulative risk. While the resulting bounds may not always be as tight as those obtained by more direct methods for specific Θ , the approach taken here does have the advantage of giving a simple, unified and general treatment to this problem, moreover, one in which no more sophisticated mathematical methods than Jensen's inequality are needed to derive the results. In the

future we hope to further explore the applications of these results to specific estimation problems, such as the “concept learning” or “pattern classification” problems examined in current machine learning and neural network research. Some initial results along these lines can be found in [46, 36] (see also [28, 43]).

There are also several other directions for further research one might pursue. Apart from general tightening of the bounds, these include treating the case of nonindependent observations, extending the results giving bounds for individual θ^* in Theorems 2 and 3 to the case where P_{θ^*} is not a distribution in Θ but is “close to” a distribution Θ , and giving a more complete characterization of the mutual information $I(\Theta^*; Y^n)$ in terms of the metric entropy properties of Θ for the infinite dimensional case, as was done for the finite dimensional case in [36].

11 Appendix

Here we give the proof of Lemma 5.

Lemma 12 *Assume $0 < \alpha < 1$ and $\lambda > 0$. Let P , R and U be any distributions on Y . Let $c_\lambda = \int dP^{1+\lambda} dU^{-\lambda}$. Let $Q = (1 - \epsilon)R + \epsilon U$ for some $\epsilon > 0$ such that $\frac{\log \log(1/\epsilon)}{\log(1/\epsilon)} \leq \lambda/2$ and $\epsilon \leq e^{-\alpha/(2(1-\alpha))}$. Then*

$$D_{KL}(P||Q) \leq \frac{2 \log \frac{1}{\epsilon}}{f_\alpha(\epsilon^2)} D_\alpha(P, R) + \frac{2\epsilon \log \frac{1}{\epsilon}}{(1-\alpha)f_\alpha(\epsilon^2)} + \epsilon^{\lambda/2} c_\lambda,$$

where

$$f_\alpha(x) = \frac{\alpha + (1-\alpha)x - x^{1-\alpha}}{1-\alpha}.$$

Proof. We use the easily verified fact that for $0 < \alpha < 1$ and $0 < x < 1$, $f_\alpha(x)$ is positive and decreasing in x . Let $Y_0 = \{y : dP(y) = 0\}$. For $y \in Y - Y_0$, let $S(y) = \frac{dQ(y)}{dP(y)}$ and $T(y) = \frac{dU(y)}{dP(y)}$. Then using Equations (3) and (4), and the definition of b_α , we have

$$D_{KL}(P||Q) = \int_{Y-Y_0} dP b_\alpha(S) f_\alpha(S) + \int_{Y_0} dQ \quad (20)$$

Consider two cases for $y \in Y - Y_0$.

1. $S(y) > \epsilon^2$ or $T(y) > \epsilon$. Here we note that since $S(y) = \frac{(1-\epsilon)dR(y) + \epsilon dU(y)}{dP(y)} \geq \epsilon \frac{dU(y)}{dP(y)} = \epsilon T(y)$, in either case $S(y) > \epsilon^2$. Hence

$$b_\alpha(S(y)) \leq b_\alpha(\epsilon^2) = \frac{\epsilon^2 + 2 \log \frac{1}{\epsilon} - 1}{f_\alpha(\epsilon^2)} \leq \frac{2 \log \frac{1}{\epsilon}}{f_\alpha(\epsilon^2)} \quad (21)$$

since b_α is decreasing.

2. $S(y) \leq \epsilon^2$ and $T(y) \leq \epsilon$. In this case

$$b_\alpha(S(y)) = \frac{S(y) + \log \frac{1}{S(y)} - 1}{f_\alpha(S(y))} \leq \frac{\log \frac{1}{S(y)}}{f_\alpha(S(y))} \leq \frac{\log \frac{1}{\epsilon} + \log \frac{1}{T(y)}}{f_\alpha(S(y))} \leq \frac{\log \frac{1}{\epsilon}}{f_\alpha(\epsilon^2)} + \frac{\log \frac{1}{T(y)}}{f_\alpha(S(y))}, \quad (22)$$

where in the last inequality we use the fact that $S(y) \leq \epsilon^2$ and $f_\alpha(x)$ is decreasing in x for $0 < x < 1$, and in the previous inequality we use the fact that $S(y) \geq \epsilon T(y)$ and that $\log(x)$ is increasing.

Let

$$W(\epsilon) = \int_{y: S(y) \leq \epsilon^2 \text{ and } T(y) \leq \epsilon} dP \log \frac{1}{T}.$$

From (20), (21), and (22) it follows that

$$D_{KL}(P||Q) \leq \frac{2 \log \frac{1}{\epsilon}}{f_\alpha(\epsilon^2)} \int_{Y-Y_0} dP f_\alpha(S) + \int_{Y_0} dQ + W(\epsilon) \leq \frac{2 \log \frac{1}{\epsilon}}{f_\alpha(\epsilon^2)} D_\alpha(P, Q) + W(\epsilon), \quad (23)$$

since $D_\alpha(P, Q) = \int_{Y-Y_0} dP f_\alpha(S) + \int_{Y_0} dQ$ and $\frac{2 \log \frac{1}{\epsilon}}{f_\alpha(\epsilon^2)} \geq 1$.

Note now that

$$\begin{aligned} D_\alpha(P, Q) &= \frac{1}{1-\alpha} \left(1 - \int (dP)^\alpha ((1-\epsilon)dR + \epsilon dU)^{1-\alpha} \right) \\ &\leq \frac{1}{1-\alpha} \left(1 - \int (dP)^\alpha ((1-\epsilon)dR)^{1-\alpha} \right) \\ &\leq \frac{1}{1-\alpha} \left(1 - \int (dP)^\alpha (dR)^{1-\alpha} \right) + \frac{1}{1-\alpha} \left(1 - (1-\epsilon)^{1-\alpha} \right) \\ &= D_\alpha(P, R) + \frac{1}{1-\alpha} \left(1 - (1-\epsilon)^{1-\alpha} \right) \\ &\leq D_\alpha(P, R) + \frac{\epsilon}{1-\alpha} \end{aligned}$$

Hence

$$D_{KL}(P||Q) \leq \frac{2 \log \frac{1}{\epsilon}}{f_\alpha(\epsilon^2)} D_\alpha(P, R) + \frac{2\epsilon \log \frac{1}{\epsilon}}{(1-\alpha)f_\alpha(\epsilon^2)} + W(\epsilon) \quad (24)$$

Finally, note that $T(y) \leq \epsilon$ implies that $\left(\frac{\epsilon}{T(y)}\right)^{\lambda/2} \geq 1$ and $\frac{\log \log(1/\epsilon)}{\log(1/\epsilon)} \leq \lambda/2$ implies that $\log \frac{1}{\gamma} \leq \left(\frac{1}{\gamma}\right)^{\lambda/2}$ for all $\gamma \leq \epsilon$. Hence when $\frac{\log \log(1/\epsilon)}{\log(1/\epsilon)} \leq \lambda/2$,

$$\begin{aligned} W(\epsilon) &= \int_{y: S(y) \leq \epsilon^2 \text{ and } T(y) \leq \epsilon} dP \log \frac{1}{T} \\ &\leq \int_{y: S(y) \leq \epsilon^2 \text{ and } T(y) \leq \epsilon} dP \left(\frac{\epsilon}{T}\right)^{\lambda/2} \log \frac{1}{T} \\ &\leq \epsilon^{\lambda/2} \int_{y: S(y) \leq \epsilon^2 \text{ and } T(y) \leq \epsilon} dP \left(\frac{1}{T}\right)^\lambda \end{aligned}$$

$$\begin{aligned}
&\leq \epsilon^{\lambda/2} \int_{y: S(y) \leq \epsilon^2 \text{ and } T(y) \leq \epsilon} dP^{1+\lambda} dU^{-\lambda} \\
&\leq \epsilon^{\lambda/2} \int dP^{1+\lambda} dU^{-\lambda} \\
&= \epsilon^{\lambda/2} c_\lambda
\end{aligned}$$

The result follows then from Inequality (24). \square

Acknowledgements

The authors would like to thank Andrew Barron for inspiring them to work on these problems and for many helpful discussions of these ideas. We also thank Shun-ichi Amari, Meir Feder, Yoav Freund, Michael Kearns, Sebastian Seung, Tom Cover, Bin Yu and Lucien Le Cam for helpful conversations, and Laszlo Györfi and an anonymous referee for their comments on an earlier version of this paper.

References

- [1] A. Barron and L. Györfi and E. van der Meulen. Distribution estimation consistent in total variation and in two types of information divergence. *IEEE Transactions on Information Theory*, 38:1437–1454, 1992.
- [2] S. Amari. Differential geometry of curved exponential families—curvatures and information loss. *Annals of Statistics*, 10:357–385, 1982.
- [3] S. Amari and N. Murata. Statistical theory of learning curves under entropic loss. *Neural Computation*, 5:140–153, 1993.
- [4] A. Barron. The strong ergodic theorem for densities: generalized Shannon-McMillan-Breiman theorem. *The Annals of Probability*, 13:1292–1303, 1985.
- [5] A. Barron. In T. M. Cover and B. Gopinath, editors, *Open Problems in Communication and Computation*, chapter 3.20. Are Bayes rules consistent in information?, pages 85–91. 1987.
- [6] A. Barron. The exponential convergence of posterior probabilities with implications for Bayes estimators of density functions. Technical Report 7, Dept. of Statistics, U. Ill. Urbana-Champaign, 1987.
- [7] A. Barron, B. Clarke, and D. Haussler. Information bounds for the risk of bayesian predictions and the redundancy of universal codes. In *Proc. International Symposium on Information Theory*.
- [8] A. Barron, B. Clarke, and D. Haussler. Information bounds for the risk of bayesian predictions and the redundancy of universal codes. In *Proc. International Symposium on Information Theory*, Jan. 1993.

- [9] A. Barron and T. Cover. A bound on the financial value of information. *IEEE Trans. on Information Theory*, 34:1097–1100, 1988.
- [10] A. Barron and Y. Yang. Information theoretic lower bounds on convergence rates of nonparametric estimators, 1995. unpublished manuscript.
- [11] L. Birgé. Approximation dans les espaces métriques et théorie de l'estimation. *Zeitschrift fuer Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 65:181–237, 1983.
- [12] L. Birgé. On estimating a density using Hellinger distance and some other strange facts. *Probability theory and related fields*, 71:271–291, 1986.
- [13] L. Birgé and P. Massart. Rates of convergence for minimum contrast estimators. *Probability Theory and Related Fields*, 97:113–150, 1993.
- [14] L. L. Cam. An extension of Wald's theory of statistical decision functions. *Annals of Mathematical Statistics*, 26:69–81, 1955.
- [15] L. L. Cam. Convergence of estimates under dimensionality restrictions. *Annals of Statistics*, 1:38–53, 1973.
- [16] R. H. Cameron and W. T. Martin. Transformation of wiener integrals under translations. *Ann. Math.*, 45:386–396, 1944.
- [17] B. Clarke. *Asymptotic cumulative risk and Bayes risk under entropy loss with applications*. PhD thesis, Dept. of Statistics, University of Ill., 1989.
- [18] B. Clarke and A. Barron. Information-theoretic asymptotics of Bayes methods. *IEEE Transactions on Information Theory*, 36(3):453–471, 1990.
- [19] B. Clarke and A. Barron. Jefferys' prior is asymptotically least favorable under entropy risk. *J. Statistical Planning and Inference*, 41:37–60, 1994.
- [20] G. F. Clements. Entropy of several sets of real-valued functions. *Pacific J. Math.*, 13:1085–1095, 1963.
- [21] T. Cover and J. Thomas. *Elements of Information Theory*. Wiley, 1991.
- [22] I. Csiszàr. On the topological properties of f -divergences. *Studia Scientiarum Math. Hungarica*, 2:329–339, 1967.
- [23] L. Davisson and A. Leon-Garcia. A source matching approach to finding minimax codes. *IEEE transactions on information theory*, IT-26:166–174, 1980.
- [24] L. Devroye and L. Györfi. *Nonparametric density estimation, the L_1 view*. Wiley, 1985.
- [25] P. Diaconis and D. Freedman. On the consistency of Bayes estimates. *Ann. Statist.*, 14:1–26, 1986.

- [26] R. M. Dudley. A course on empirical processes. *Lecture Notes in Mathematics*, 1097:2–142, 1984.
- [27] S. Y. Efroimovich. Information contained in a sequence of observations. *Problems in Information Transmission*, 15:178–189, 1980.
- [28] M. Feder, Y. Freund, and Y. Mansour. Optimal universal learning and prediction of probabilistic concepts. In *Proc. of IEEE Information Theory Conference*, page 233. IEEE, 1995.
- [29] R. Gallager. Source coding with side information and universal coding. Technical Report LIDS-P-937, MIT Laboratory for Information and Decision Systems, 1979.
- [30] J. Ghosh, S. Ghosal, and T. Samanta. Statistical decision theory and related topics v. In S. Gupta and J. O. Berger, editors, *Stability and Convergence of the Posterior in Non-Regular Problems*. Springer Verlag.
- [31] E. Giné and J. Zinn. Some limit theorems for empirical processes. *Annals of Probability*, 12:929–989, 1984.
- [32] R. Hasminskii and I. Ibragimov. On density estimation in the view of Kolmogorov’s ideas in approximation theory. *Annals of statistics*, 18:999–1010, 1990.
- [33] D. Haussler. A general minimax result for relative entropy. *IEEE Trans. Info Th.*, 1997. to appear.
- [34] D. Haussler and A. Barron. How well do Bayes methods work for on-line prediction of $\{+1, -1\}$ values? In *Proceedings of the Third NEC Symposium on Computation and Cognition*. SIAM, 1992.
- [35] D. Haussler, M. Kearns, and R. E. Schapire. Bounds on the sample complexity of Bayesian learning using information theory and the VC dimension. *Machine Learning*, 14(1):83–113, 1994.
- [36] D. Haussler and M. Opper. General bounds on the mutual information between a parameter and n conditionally independent observations. In *Proceedings of the Seventh Annual ACM Workshop on Computational Learning Theory*, 1995.
- [37] I. Ibragimov and R. Hasminskii. On the information in a sample about a parameter. In *Second Int. Symp. on Information Theory*, pages 295–309, 1972.
- [38] A. J. Izenman. Recent developments in nonparametric density estimation. *JASA*, 86(413):205–224, 1991.
- [39] A. N. Kolmogorov and V. M. Tihomirov. ϵ -entropy and ϵ -capacity of sets in functional spaces. *Amer. Math. Soc. Translations (Ser. 2)*, 17:277–364, 1961.
- [40] F. Komaki. On asymptotic properties of predictive distributions. Technical Report METR 94-21, U. Yokyo, Math. and Eng. Physics, 1994.

- [41] L. Györfi and I. Páli and E. van der Meulen. There is no universal source code for an infinite alphabet. *IEEE Transactions on Information Theory*, 40:267–271, 1994.
- [42] L. LeCam. *Asymptotic methods in statistical decision theory*. Springer, 1986.
- [43] R. Meir and N. Merhav. On the stochastic complexity of learning realizable and unrealizable rules. *Machine Learning*, 19(3):241–, 1995.
- [44] N. Merhav and M. Feder. A strong version of the redundancy-capacity theorem of universal coding. *IEEE Trans. Info Th.*, 41(3):714–, 1995.
- [45] M. Opper and D. Haussler. Calculation of the learning curve of Bayes optimal classification algorithm for learning a perceptron with noise. In *Proc. 4th Annu. Workshop on Comput. Learning Theory*, pages 75–87, San Mateo, CA, 1991. Morgan Kaufmann.
- [46] M. Opper and D. Haussler. Bounds for predictive errors in the statistical mechanics of in supervised learning. *Physical Review Letters*, 75(20):3772–3775, 1995.
- [47] M. S. Pinsker. *Information and Information Stability of Random Variables and Processes (Transl.)*. Holden Day, 1964.
- [48] D. Pollard. *Empirical Processes: Theory and Applications*, volume 2 of *NSF-CBMS Regional Conference Series in Probability and Statistics*. Institute of Math. Stat. and Am. Stat. Assoc., 1990.
- [49] A. Renyi. On measures of entropy and information. In L. C. et.al., editor, *Fourth Berkeley Sym. on Math., Stat. and Prob.*, pages 547–561. 1960.
- [50] A. Renyi. On the amount of information concerning an unknown parameter in a sequence of observations. *Publ. Math. Inst. Hungar. Acad. Sci.*, 9:617–625, 1964.
- [51] J. Rissanen. Stochastic complexity and modeling. *The Annals of Statistics*, 14(3):1080–1100, 1986.
- [52] J. Rissanen, T. Speed, and B. Yu. Density estimation by stochastic complexity. *IEEE Trans. Info. Th.*, 38:315–323, 1992.
- [53] K. Symanzik. Proof and refinements of an inequality of Feynman. *J.Math. Phys.*, 6:1155–, 1965.
- [54] S. van deGeer. Hellinger-consistency of certain nonparametric maximum likelihood estimators. *Annals of Statistics*, 21:14–44, 1993.
- [55] W. Wong and X. Shen. Probability inequalities for likelihood ratios and convergence rates for sieve MLE's. *Annals of Statistics*, 23(2):339–362, 1995.
- [56] K. Yamanishi. A learning criterion for stochastic rules. *Machine Learning*, 1992. Special Issue on the Proceedings of the 3rd Workshop on Computational Learning Theory.

- [57] K. Yamanishi. A decision-theoretic extension of stochastic complexity and its applications to learning. Unpublished manuscript, 1995.
- [58] K. Yamanishi. A loss bound model for on-line stochastic prediction algorithms. *Information and Computation*, 119:39–54, 1995.
- [59] B. Yu. Lower bounds on expected redundancy for nonparametric classes. *IEEE Trans. Info. Th.*, 42(1), 1996.
- [60] H. Zhu and R. Rohwer. Information geometric measurements of generalization. Technical Report NCRG 4350, Aston U., England, Neural computing research group, 1995.