

---

# MDL Convergence Speed for Bernoulli Sequences\*

---

**Jan Poland** and **Marcus Hutter**

IDSIA, Galleria 2 CH-6928 Manno (Lugano), Switzerland<sup>†</sup>  
{jan,marcus}@idsia.ch <http://www.idsia.ch>

22 February 2006

## Abstract

The Minimum Description Length principle for online sequence estimation/prediction in a proper learning setup is studied. If the underlying model class is discrete, then the total expected square loss is a particularly interesting performance measure: (a) this quantity is finitely bounded, implying convergence with probability one, and (b) it additionally specifies the convergence speed. For MDL, in general one can only have loss bounds which are finite but exponentially larger than those for Bayes mixtures. We show that this is even the case if the model class contains only Bernoulli distributions. We derive a new upper bound on the prediction error for countable Bernoulli classes. This implies a small bound (comparable to the one for Bayes mixtures) for certain important model classes. We discuss the application to Machine Learning tasks such as classification and hypothesis testing, and generalization to countable classes of i.i.d. models.

## Keywords

MDL, Minimum Description Length, Convergence Rate, Prediction, Bernoulli, Discrete Model Class.

---

\*A shorter version of this paper [PH04b] appeared in ALT 2004.

<sup>†</sup>This work was supported by SNF grant 2100-67712.02.

# 1 Introduction

“Bayes mixture”, “Solomonoff induction”, “marginalization”, all these terms refer to a central induction principle: Obtain a predictive distribution by integrating the product of prior and evidence over the model class. In many cases however, the Bayes mixture is computationally infeasible, and even a sophisticated approximation is expensive. The MDL or MAP (maximum a posteriori) estimator is both a common approximation for the Bayes mixture and interesting for its own sake: Use the model with the largest product of prior and evidence. (In practice, the MDL estimator is usually being approximated too, since only a local maximum is determined.)

How good are the predictions by Bayes mixtures and MDL? This question has attracted much attention. In many cases, an important quality measure is the *total* or cumulative *expected loss* of a predictor. In particular the square loss is often considered. Assume that the outcome space is finite, and the model class is continuously parameterized. Then for Bayes mixture prediction, the cumulative expected square loss is usually small but unbounded, growing with  $\ln n$ , where  $n$  is the sample size [CB90, Hut03b]. This corresponds to an *instantaneous* loss bound of  $\frac{1}{n}$ . For the MDL predictor, the losses behave similarly [Ris96, BRY98] under appropriate conditions, in particular with a specific prior. (Note that in order to do MDL for continuous model classes, one needs to *discretize* the parameter space, see also [BC91].)

On the other hand, if the model class is discrete, then Solomonoff’s theorem [Sol78, Hut01] bounds the cumulative expected square loss for the Bayes mixture predictions finitely, namely by  $\ln w_\mu^{-1}$ , where  $w_\mu$  is the prior weight of the “true” model  $\mu$ . The only necessary assumption is that the true distribution  $\mu$  is contained in the model class, i.e. that we are dealing with *proper learning*. It has been demonstrated [GL04], that for both Bayes mixture and MDL, the proper learning assumption can be essential: If it is violated, then learning may fail very badly.

For MDL predictions in the proper learning case, it has been shown [PH04a] that a bound of  $w_\mu^{-1}$  holds. This bound is exponentially larger than the Solomonoff bound, and it is sharp in general. A finite bound on the total expected square loss is particularly interesting:

1. It implies convergence of the predictive to the true probabilities with probability one. In contrast, an instantaneous loss bound of  $\frac{1}{n}$  implies only convergence in probability.
2. Additionally, it gives a *convergence speed*, in the sense that errors of a certain magnitude cannot occur too often.

So for both, Bayes mixtures and MDL, convergence with probability one holds, while the convergence speed is exponentially worse for MDL compared to the Bayes mixture. (We avoid the term “convergence rate” here, since the order of convergence is identical in both cases. It is e.g.  $o(1/n)$  if we additionally assume that the error is monotonically decreasing, which is not necessarily true in general).

It is therefore natural to ask if there are model classes where the cumulative loss of MDL is comparable to that of Bayes mixture predictions. In the present work, we concentrate on the simplest possible stochastic case, namely discrete Bernoulli classes. (Note that then the MDL “predictor” just becomes an estimator, in that it estimates the true parameter and directly uses that for prediction. Nevertheless, for consistency of terminology, we keep the term predictor.) It might be surprising to discover that in general the cumulative loss is still exponential. On the other hand, we will give mild conditions on the prior guaranteeing a small bound. Moreover, it is well-known that the instantaneous square loss of the Maximum Likelihood estimator decays as  $\frac{1}{n}$  in the Bernoulli case. The same holds for MDL, as we will see. (If convergence speed is measured in terms of instantaneous losses, then much more general statements are possible [Li99, Zha04], this is briefly discussed in Section 4.)

A particular motivation to consider discrete model classes arises in Algorithmic Information Theory. From a computational point of view, the largest relevant model class is the class of all computable models on some fixed universal Turing machine, precisely prefix machine [LV97]. Thus each model corresponds to a program, and there are countably many programs. Moreover, the models are stochastic, precisely they are *semimeasures* on strings (programs need not halt, otherwise the models were even measures). Each model has a natural description length, namely the length of the corresponding program. If we agree that programs are binary strings, then a prior is defined by two to the negative description length. By the Kraft inequality, the priors sum up to at most one.

Also the Bernoulli case can be studied in the view of Algorithmic Information Theory. We call this the *universal setup*: Given a universal Turing machine, the related class of Bernoulli distributions is isomorphic to the countable set of computable reals in  $[0, 1]$ . The description length  $Kw(\vartheta)$  of a parameter  $\vartheta \in [0, 1]$  is then given by the length of its shortest program. A prior weight may then be defined by  $2^{-Kw(\vartheta)}$ . (If a string  $x = x_1x_2\dots x_{t-1}$  is generated by a Bernoulli distribution with computable parameter  $\vartheta_0 \in [0, 1]$ , then with high probability the two-part complexity of  $x$  with respect to the Bernoulli class does not exceed its algorithmic complexity by more than a constant, as shown by Vovk [Vov97]. That is, the two-part complexity with respect to the Bernoulli class *is* the shortest description, save for an additive constant.)

Many Machine Learning tasks are or can be reduced to sequence prediction tasks. An important example is classification. The task of classifying a new instance  $z_n$  after having seen (instance,class) pairs  $(z_1, c_1), \dots, (z_{n-1}, c_{n-1})$  can be phrased as to predict the continuation of the sequence  $z_1c_1\dots z_{n-1}c_{n-1}z_n$ . Typically the (instance,class) pairs are i.i.d. Cumulative loss bounds for prediction usually generalize to prediction *conditionalized* to some inputs [PH05]. Then we can solve classification problems in the standard form. It is not obvious if and how the proofs in this paper can be conditionalized.

Our main tool for obtaining results is the Kullback-Leibler divergence. Lemmata for this quantity are stated in Section 2. Section 3 shows that the exponential error

bound obtained in [PH04a] is sharp in general. In Section 4, we give an upper bound on the instantaneous and the cumulative losses. The latter bound is small e.g. under certain conditions on the distribution of the weights, this is the subject of Section 5. Section 6 treats the universal setup. Finally, in Section 7 we discuss the results and give conclusions.

## 2 Kullback-Leibler Divergence

Let  $\mathbb{B} = \{0, 1\}$  and consider finite strings  $x \in \mathbb{B}^*$  as well as infinite sequences  $x_{<\infty} \in \mathbb{B}^\infty$ , with the first  $n$  bits denoted by  $x_{1:n}$ . If we know that  $x$  is generated by an i.i.d random variable, then  $P(x_i = 1) = \vartheta_0$  for all  $1 \leq i \leq \ell(x)$  where  $\ell(x)$  is the length of  $x$ . Then  $x$  is called a Bernoulli sequence, and  $\vartheta_0 \in \Theta \subset [0, 1]$  the *true parameter*. In the following we will consider only countable  $\Theta$ , e.g. the set of all computable numbers in  $[0, 1]$ .

Associated with each  $\vartheta \in \Theta$ , there is a *complexity* or description length  $Kw(\vartheta)$  and a *weight* or (semi)probability  $w_\vartheta = 2^{-Kw(\vartheta)}$ . The complexity will often but need not be a natural number. Typically, one assumes that the weights sum up to at most one,  $\sum_{\vartheta \in \Theta} w_\vartheta \leq 1$ . Then, by the Kraft inequality, for all  $\vartheta \in \Theta$  there exists a prefix-code of length  $Kw(\vartheta)$ . Because of this correspondence, it is only a matter of convenience whether results are developed in terms of description lengths or probabilities. We will choose the former way. We won't even need the condition  $\sum_{\vartheta} w_\vartheta \leq 1$  for most of the following results. This only means that  $Kw$  cannot be interpreted as a prefix code length, but does not cause other problems.

Given a set of distributions  $\Theta \subset [0, 1]$ , complexities  $(Kw(\vartheta))_{\vartheta \in \Theta}$ , a true distribution  $\vartheta_0 \in \Theta$ , and some observed string  $x \in \mathbb{B}^*$ , we define an *MDL estimator*<sup>1</sup>:

$$\vartheta^x = \arg \max_{\vartheta \in \Theta} \{w_\vartheta P(x|\vartheta)\}.$$

Here,  $P(x|\vartheta)$  is the probability of observing  $x$  if  $\vartheta$  is the true parameter. Clearly,  $P(x|\vartheta) = \vartheta^{\mathbb{1}(x)}(1 - \vartheta)^{\ell(x) - \mathbb{1}(x)}$ , where  $\mathbb{1}(x)$  is the number of ones in  $x$ . Hence  $P(x|\vartheta)$  depends only on  $\ell(x)$  and  $\mathbb{1}(x)$ . We therefore see

$$\begin{aligned} \vartheta^x &= \vartheta^{(\alpha, n)} = \arg \max_{\vartheta \in \Theta} \{w_\vartheta (\vartheta^\alpha (1 - \vartheta)^{1-\alpha})^n\} \\ &= \arg \min_{\vartheta \in \Theta} \{n \cdot D(\alpha \parallel \vartheta) + Kw(\vartheta) \cdot \ln 2\}, \end{aligned} \tag{1}$$

where  $n = \ell(x)$  and  $\alpha := \frac{\mathbb{1}(x)}{\ell(x)}$  is the *observed fraction* of ones and

$$D(\alpha \parallel \vartheta) = \alpha \ln \frac{\alpha}{\vartheta} + (1 - \alpha) \ln \frac{1 - \alpha}{1 - \vartheta}$$

---

<sup>1</sup>Precisely, we define a MAP (maximum a posteriori) estimator. For two reasons, our definition might not be considered as MDL in the strict sense. First, MDL is often associated with a specific prior, while we admit arbitrary priors. Second and more importantly, when coding some data  $x$ , one can exploit the fact that once the parameter  $\vartheta^x$  is specified, only data which leads to this  $\vartheta^x$  needs to be considered. This allows for a description shorter than  $Kw(\vartheta^x)$ . Nevertheless, the *construction principle* is commonly termed MDL, compare e.g. the “ideal MDL” in [VL00].

is the Kullback-Leibler divergence. The second line of (1) also explains the name MDL, since we choose the  $\vartheta$  which minimizes the joint description of model  $\vartheta$  and the data  $x$  given the model.

We also define the *extended Kullback-Leibler divergence*

$$D^\alpha(\vartheta\|\tilde{\vartheta}) = \alpha \ln \frac{\vartheta}{\tilde{\vartheta}} + (1 - \alpha) \ln \frac{1 - \vartheta}{1 - \tilde{\vartheta}} = D(\alpha\|\tilde{\vartheta}) - D(\alpha\|\vartheta). \quad (2)$$

It is easy to see that  $D^\alpha(\vartheta\|\tilde{\vartheta})$  is linear in  $\alpha$ ,  $D^\vartheta(\vartheta\|\tilde{\vartheta}) = D(\vartheta\|\tilde{\vartheta})$  and  $D^{\tilde{\vartheta}}(\vartheta\|\tilde{\vartheta}) = -D(\tilde{\vartheta}\|\vartheta)$ , and  $\frac{d}{d\alpha}D^\alpha(\vartheta\|\tilde{\vartheta}) > 0$  iff  $\vartheta > \tilde{\vartheta}$ . Note that  $D^\alpha(\vartheta\|\tilde{\vartheta})$  may be also defined for the general i.i.d. case, i.e. if the alphabet has more than two symbols.

Let  $\vartheta, \tilde{\vartheta} \in \Theta$  be two parameters, then it follows from (1) that in the process of choosing the MDL estimator,  $\vartheta$  is being preferred to  $\tilde{\vartheta}$  iff

$$nD^\alpha(\vartheta\|\tilde{\vartheta}) \geq \ln 2 \cdot (Kw(\vartheta) - Kw(\tilde{\vartheta})) \quad (3)$$

with  $n$  and  $\alpha$  as before. We also say that then  $\vartheta$  *beats*  $\tilde{\vartheta}$ . It is immediate that for increasing  $n$  the influence of the complexities on the selection of the maximizing element decreases. We are now interested in the *total expected square prediction error* (or cumulative square loss) of the MDL estimator

$$\sum_{n=1}^{\infty} \mathbf{E}(\vartheta^{x_{1:n}} - \vartheta_0)^2.$$

In terms of [PH04a], this is the *static MDL prediction* loss, which means that a predictor/estimator  $\vartheta^x$  is chosen according to the current observation  $x$ . (As already mentioned, the terms predictor and estimator coincide for static MDL and Bernoulli classes.) The *dynamic* method on the other hand would consider both possible continuations  $x_0$  and  $x_1$  and predict according to  $\vartheta^{x_0}$  and  $\vartheta^{x_1}$ . In the following, we concentrate on static predictions. They are also preferred in practice, since computing only one model is more efficient.

Let  $A_n = \{\frac{k}{n} : 0 \leq k \leq n\}$ . Given the true parameter  $\vartheta_0$  and some  $n \in \mathbb{N}$ , the *expectation* of a function  $f^{(n)} : \{0, \dots, n\} \rightarrow \mathbb{R}$  is given by

$$\mathbf{E}f^{(n)} = \sum_{\alpha \in A_n} p(\alpha|n)f(\alpha n), \text{ where } p(\alpha|n) = \binom{n}{k} \left( \vartheta_0^\alpha (1 - \vartheta_0)^{1-\alpha} \right)^n. \quad (4)$$

(Note that the probability  $p(\alpha|n)$  depends on  $\vartheta_0$ , which we do not make explicit in our notation.) Therefore,

$$\sum_{n=1}^{\infty} \mathbf{E}(\vartheta^{x_{1:n}} - \vartheta_0)^2 = \sum_{n=1}^{\infty} \sum_{\alpha \in A_n} p(\alpha|n)(\vartheta^{(\alpha,n)} - \vartheta_0)^2, \quad (5)$$

Denote the relation  $f = O(g)$  by  $f \overset{\times}{\leq} g$ . Analogously define “ $\overset{\times}{\geq}$ ” and “ $\overset{\times}{=}$ ”. From [PH04a, Corollary 12], we immediately obtain the following result.

**Theorem 1** *The cumulative loss bound  $\sum_n \mathbf{E}(\vartheta^{x_{1:n}} - \vartheta_0)^2 \leq 2^{Kw(\vartheta_0)}$  holds.*

This is the “slow” convergence result mentioned in the introduction. In contrast, for a Bayes mixture, the total expected error is bounded by  $Kw(\vartheta_0)$  rather than  $2^{Kw(\vartheta_0)}$  (see [Sol78] or [Hut01, Th.1]). An upper bound on  $\sum_n \mathbf{E}(\vartheta^{x_{1:n}} - \vartheta_0)^2$  is termed as *convergence in mean sum* and implies convergence  $\vartheta^{x_{1:n}} \rightarrow \vartheta_0$  with probability 1 (since otherwise the sum would be infinite).

We now establish relations between the Kullback-Leibler divergence and the quadratic distance. We call bounds of this type *entropy inequalities*.

**Lemma 2** *Let  $\vartheta, \tilde{\vartheta} \in (0, 1)$ . Let  $\vartheta^* = \arg \min\{|\vartheta - \frac{1}{2}|, |\tilde{\vartheta} - \frac{1}{2}|\}$ , i.e.  $\vartheta^*$  is the element from  $\{\vartheta, \tilde{\vartheta}\}$  which is closer to  $\frac{1}{2}$ . Then the following assertions hold.*

- (i)  $D(\vartheta \parallel \tilde{\vartheta}) \geq 2 \cdot (\vartheta - \tilde{\vartheta})^2 \quad \forall \vartheta, \tilde{\vartheta} \in (0, 1),$
- (ii)  $D(\vartheta \parallel \tilde{\vartheta}) \leq \frac{8}{3}(\vartheta - \tilde{\vartheta})^2 \quad \text{if } \vartheta, \tilde{\vartheta} \in [\frac{1}{4}, \frac{3}{4}],$
- (iii)  $D(\vartheta \parallel \tilde{\vartheta}) \geq \frac{(\vartheta - \tilde{\vartheta})^2}{2\vartheta^*(1 - \vartheta^*)} \quad \text{if } \vartheta, \tilde{\vartheta} \leq \frac{1}{2},$
- (iv)  $D(\vartheta \parallel \tilde{\vartheta}) \leq \frac{3(\vartheta - \tilde{\vartheta})^2}{2\vartheta^*(1 - \vartheta^*)} \quad \text{if } \vartheta \leq \frac{1}{4} \text{ and } \tilde{\vartheta} \in [\frac{\vartheta}{3}, 3\vartheta],$
- (v)  $D(\tilde{\vartheta} \parallel \vartheta) \geq \tilde{\vartheta}(\ln \tilde{\vartheta} - \ln \vartheta - 1) \quad \forall \vartheta, \tilde{\vartheta} \in (0, 1),$
- (vi)  $D(\vartheta \parallel \tilde{\vartheta}) \leq \frac{1}{2}\tilde{\vartheta} \quad \text{if } \vartheta \leq \tilde{\vartheta} \leq \frac{1}{2},$
- (vii)  $D(\vartheta \parallel \vartheta \cdot 2^{-j}) \leq j \cdot \vartheta \quad \text{if } \vartheta \leq \frac{1}{2} \text{ and } j \geq 1,$
- (viii)  $D(\vartheta \parallel 1 - 2^{-j}) \leq j \quad \text{if } \vartheta \leq \frac{1}{2} \text{ and } j \geq 1.$

*Statements (iii) – (viii) have symmetric counterparts for  $\vartheta \geq \frac{1}{2}$ .*

The first two statements give upper and lower bounds for the Kullback-Leibler divergence in terms of the quadratic distance. They express the fact that the Kullback-Leibler divergence is locally quadratic. So do the next two statements, they will be applied in particular if  $\vartheta$  is located close to the boundary of  $[0, 1]$ . Statements (v) and (vi) give bounds in terms of the absolute distance, i.e. “linear” instead of quadratic. They are mainly used if  $\tilde{\vartheta}$  is relatively far from  $\vartheta$ . Note that in (v), the position of  $\vartheta$  and  $\tilde{\vartheta}$  are inverted. The last two inequalities finally describe the behavior of the Kullback-Leibler divergence as its second argument tends to the boundary of  $[0, 1]$ . Observe that this is logarithmic in the inverse distance to the boundary.

**Proof.** (i) This is standard, see e.g. [LV97]. It is shown similarly as (iii).

(ii) Let  $f(\eta) = D(\vartheta \parallel \eta) - \frac{8}{3}(\eta - \vartheta)^2$ , then we show  $f(\eta) \leq 0$  for  $\eta \in [\frac{1}{4}, \frac{3}{4}]$ . We have that  $f(\vartheta) = 0$  and

$$f'(\eta) = \frac{\eta - \vartheta}{\eta(1 - \eta)} - \frac{16}{3}(\eta - \vartheta).$$

This difference is nonnegative if and only  $\eta - \vartheta \leq 0$  since  $\eta(1 - \eta) \geq \frac{3}{16}$ . This implies  $f(\eta) \leq 0$ .

(iii) Consider the function

$$f(\eta) = D(\vartheta\|\eta) - \frac{(\vartheta - \eta)^2}{2 \max\{\vartheta, \eta\}(1 - \max\{\vartheta, \eta\})}.$$

We have to show that  $f(\eta) \geq 0$  for all  $\eta \in (0, \frac{1}{2}]$ . It is obvious that  $f(\vartheta) = 0$ . For  $\eta \leq \vartheta$ ,

$$f'(\eta) = \frac{\eta - \vartheta}{\eta(1 - \eta)} - \frac{\eta - \vartheta}{\vartheta(1 - \vartheta)} \leq 0$$

holds since  $\eta - \vartheta \leq 0$  and  $\vartheta(1 - \vartheta) \geq \eta(1 - \eta)$ . Thus,  $f(\eta) \geq 0$  must be valid for  $\eta \leq \vartheta$ . On the other hand if  $\eta \geq \vartheta$ , then

$$f'(\eta) = \frac{\eta - \vartheta}{\eta(1 - \eta)} - \left[ \frac{\eta - \vartheta}{\eta(1 - \eta)} - \frac{(\eta - \vartheta)^2(1 - 2\eta)}{2\eta^2(1 - \eta)^2} \right] \geq 0$$

is true. Thus  $f(\eta) \geq 0$  holds in this case, too.

(iv) We show that

$$f(\eta) = D(\vartheta\|\eta) - \frac{3(\vartheta - \eta)^2}{2 \max\{\vartheta, \eta\}(1 - \max\{\vartheta, \eta\})} \leq 0$$

for  $\eta \in [\frac{\vartheta}{3}, 3\vartheta]$ . If  $\eta \leq \vartheta$ , then

$$f'(\eta) = \frac{\eta - \vartheta}{\eta(1 - \eta)} - \frac{3(\eta - \vartheta)}{\vartheta(1 - \vartheta)} \geq 0$$

since  $3\eta(1 - \eta) \geq \vartheta(1 - \eta) \geq \vartheta(1 - \vartheta)$ . If  $\eta \geq \vartheta$ , then

$$f'(\eta) = \frac{\eta - \vartheta}{\eta(1 - \eta)} - 3 \cdot \left[ \frac{\eta - \vartheta}{\eta(1 - \eta)} - \frac{(\eta - \vartheta)^2(1 - 2\eta)}{2\eta^2(1 - \eta)^2} \right] \leq 0$$

is equivalent to  $4\eta(1 - \eta) \geq 3(\eta - \vartheta)(1 - 2\eta)$ , which is fulfilled if  $\vartheta \leq \frac{1}{4}$  and  $\eta \leq 3\vartheta$  as an elementary computation verifies.

(v) Using  $-\ln(1-u) \leq \frac{u}{1-u}$ , one obtains

$$\begin{aligned} D(\tilde{\vartheta}\|\vartheta) &= \tilde{\vartheta} \ln \frac{\tilde{\vartheta}}{\vartheta} + (1 - \tilde{\vartheta}) \ln \frac{1 - \tilde{\vartheta}}{1 - \vartheta} \geq \tilde{\vartheta} \ln \frac{\tilde{\vartheta}}{\vartheta} + (1 - \tilde{\vartheta}) \ln(1 - \tilde{\vartheta}) \\ &\geq \tilde{\vartheta} \ln \frac{\tilde{\vartheta}}{\vartheta} - \tilde{\vartheta} = \tilde{\vartheta}(\ln \tilde{\vartheta} - \ln \vartheta - 1) \end{aligned}$$

(vi) This follows from  $D(\vartheta\|\tilde{\vartheta}) \leq -\ln(1 - \tilde{\vartheta}) \leq \frac{\tilde{\vartheta}}{1 - \tilde{\vartheta}} \leq \frac{\tilde{\vartheta}}{2}$ . The last two statements (vii) and (viii) are even easier.  $\square$

In the above entropy inequalities we have left out the extreme cases  $\vartheta, \tilde{\vartheta} \in \{0, 1\}$ . This is for simplicity and convenience only. Inequalities (i) – (iv) remain valid for  $\vartheta, \tilde{\vartheta} \in \{0, 1\}$  if the fraction  $\frac{0}{0}$  is properly defined. However, since the extreme

cases will need to be considered separately anyway, there is no requirement for the extension of the lemma. We won't need (vi) and (viii) of Lemma 2 in the sequel.

We want to point out that although we have proven Lemma 2 only for the case of binary alphabet, generalizations to arbitrary alphabet are likely to hold. In fact, (i) does hold for arbitrary alphabet, as shown in [Hut01].

It is a well-known fact that the binomial distribution may be approximated by a Gaussian. Our next goal is to establish upper and lower bounds for the binomial distribution. Again we leave out the extreme cases.

**Lemma 3** *Let  $\vartheta_0 \in (0, 1)$  be the true parameter,  $n \geq 2$  and  $1 \leq k \leq n - 1$ , and  $\alpha = \frac{k}{n}$ . Then the following assertions hold.*

$$(i) \quad p(\alpha|n) \leq \frac{1}{\sqrt{2\pi\alpha(1-\alpha)n}} \exp(-nD(\alpha||\vartheta_0)),$$

$$(ii) \quad p(\alpha|n) \geq \frac{1}{\sqrt{8\alpha(1-\alpha)n}} \exp(-nD(\alpha||\vartheta_0)).$$

The lemma gives a quantitative assertion about the Gaussian approximation to a binomial distribution. The upper bound is sharp for  $n \rightarrow \infty$  and fixed  $\alpha$ . Lemma 3 can be easily combined with Lemma 2, yielding Gaussian estimates for the Binomial distribution.

**Proof.** Stirling's formula is a well-known result from calculus. In a refined version, it states that for any  $n \geq 1$  the factorial  $n!$  can be bounded from below and above by

$$\sqrt{2\pi n} \cdot n^n \exp\left(-n + \frac{1}{12n+1}\right) \leq n! \leq \sqrt{2\pi n} \cdot n^n \exp\left(-n + \frac{1}{12n}\right).$$

Hence,

$$\begin{aligned} p(\alpha, n) &= \frac{n!}{k!(n-k)!} \vartheta_0^k (1-\vartheta_0)^{n-k} \\ &\leq \frac{\sqrt{n} \cdot n^n \exp\left(\frac{1}{12n}\right) \vartheta_0^k (1-\vartheta_0)^{n-k}}{\sqrt{2\pi k(n-k)} \cdot k^k (n-k)^{n-k} \exp\left(\frac{1}{12k+1} + \frac{1}{12(n-k)+1}\right)} \\ &= \frac{1}{\sqrt{2\pi\alpha(1-\alpha)n}} \exp\left(-n \cdot D(\alpha||\vartheta_0) + \frac{1}{12n} - \frac{1}{12k+1} - \frac{1}{12(n-k)+1}\right) \\ &\leq \frac{1}{\sqrt{2\pi\alpha(1-\alpha)n}} \exp(-nD(\alpha||\vartheta_0)). \end{aligned}$$

The last inequality is valid since  $\frac{1}{12n} - \frac{1}{12k+1} - \frac{1}{12(n-k)+1} < 0$  for all  $n$  and  $k$ , which is easily verified using elementary computations. This establishes (i).

In order to show (ii), we observe

$$\begin{aligned} p(\alpha, n) &\geq \frac{1}{\sqrt{2\pi\alpha(1-\alpha)n}} \exp\left(-n \cdot D(\alpha\|\vartheta_0) + \frac{1}{12n+1} - \frac{1}{12k} - \frac{1}{12(n-k)}\right) \\ &\geq \frac{\exp(\frac{1}{37} - \frac{1}{8})}{\sqrt{2\pi\alpha(1-\alpha)n}} \exp(-nD(\alpha\|\vartheta_0)) \quad \text{for } n \geq 3. \end{aligned}$$

Here the last inequality follows from the fact that  $\frac{1}{12n+1} - \frac{1}{12k} - \frac{1}{12(n-k)}$  is minimized for  $n = 3$  (and  $k = 1$  or  $2$ ), if we exclude  $n = 2$ , and  $\exp(\frac{1}{37} - \frac{1}{8}) \geq \sqrt{\pi}/2$ . For  $n = 2$  a direct computation establishes the lower bound.  $\square$

**Lemma 4** Let  $z \in \mathbb{R}^+$ , then

$$\begin{aligned} (i) \quad &\frac{\sqrt{\pi}}{2z^3} - \frac{1}{z\sqrt{2e}} \leq \sum_{n=1}^{\infty} \sqrt{n} \cdot \exp(-z^2n) \leq \frac{\sqrt{\pi}}{2z^3} + \frac{1}{z\sqrt{2e}} \quad \text{and} \\ (ii) \quad &\sum_{n=1}^{\infty} n^{-\frac{1}{2}} \exp(-z^2n) \leq \sqrt{\pi}/z. \end{aligned}$$

**Proof.** (i) The function  $f(u) = \sqrt{u} \exp(-z^2u)$  increases for  $u \leq \frac{1}{2z^2}$  and decreases for  $u \geq \frac{1}{2z^2}$ . Let  $N = \max\{n \in \mathbb{N} : f(n) \geq f(n-1)\}$ , then it is easy to see that

$$\begin{aligned} \sum_{n=1}^{N-1} f(n) &\leq \int_0^N f(u) du \leq \sum_{n=1}^N f(n) \quad \text{and} \\ \sum_{n=N+1}^{\infty} f(n) &\leq \int_N^{\infty} f(u) du \leq \sum_{n=N}^{\infty} f(n) \quad \text{and thus} \\ \sum_{n=1}^{\infty} f(n) - f(N) &\leq \int_0^{\infty} f(u) du \leq \sum_{n=1}^{\infty} f(n) + f(N) \end{aligned}$$

holds. Moreover,  $f$  is the derivative of the function

$$F(u) = -\frac{\sqrt{u} \exp(-z^2u)}{z^2} + \frac{1}{z^3} \int_0^{z\sqrt{u}} \exp(-v^2) dv.$$

Observe  $f(N) \leq f(\frac{1}{2z^2}) = \frac{\exp(-\frac{1}{2})}{z\sqrt{2}}$  and  $\int_0^{\infty} \exp(-v^2) dv = \frac{\sqrt{\pi}}{2}$  to obtain the assertion.

(ii) The function  $f(u) = u^{-\frac{1}{2}} \exp(-z^2u)$  decreases monotonically on  $(0, \infty)$  and is the derivative of  $F(u) = 2z^{-1} \int_0^{z\sqrt{u}} \exp(-v^2) dv$ . Therefore,

$$\sum_{n=1}^{\infty} f(n) \leq \int_0^{\infty} f(u) du = \sqrt{\pi}/z$$

holds.  $\square$

### 3 Lower Bound

We are now in the position to prove that even for Bernoulli classes the upper bound from Theorem 1 is sharp in general.

**Proposition 5** *Let  $\vartheta_0 = \frac{1}{2}$  be the true parameter generating sequences of fair coin flips. Assume  $\Theta = \{\vartheta_0, \vartheta_1, \dots, \vartheta_{2^N-1}\}$  where  $\vartheta_k = \frac{1}{2} + 2^{-k-1}$  for  $k \geq 1$ . Let all complexities be equal, i.e.  $Kw(\vartheta_0) = \dots = Kw(\vartheta_{2^N-1}) = N$ . Then*

$$\sum_{n=1}^{\infty} \mathbf{E}(\vartheta_0 - \vartheta^x)^2 \geq \frac{1}{84}(2^N - 5) \stackrel{\times}{\approx} 2^{Kw(\vartheta_0)}.$$

**Proof.** Recall that  $\vartheta^x = \vartheta^{(\alpha, n)}$  the maximizing element for some observed sequence  $x$  only depends on the length  $n$  and the observed fraction of ones  $\alpha$ . In order to obtain an estimate for the total prediction error  $\sum_n \mathbf{E}(\vartheta_0 - \vartheta^x)^2$ , partition the interval  $[0, 1]$  into  $2^N$  disjoint intervals  $I_k$ , such that  $\bigcup_{k=0}^{2^N-1} I_k = [0, 1]$ . Then consider the contributions for the observed fraction  $\alpha$  falling in  $I_k$  separately:

$$C(k) = \sum_{n=1}^{\infty} \sum_{\alpha \in A_n \cap I_k} p(\alpha|n)(\vartheta^{(\alpha, n)} - \vartheta_0)^2 \quad (6)$$

(compare (4)). Clearly,  $\sum_n \mathbf{E}(\vartheta_0 - \vartheta^x)^2 = \sum_k C(k)$  holds. We define the partitioning ( $I_k$ ) as  $I_0 = [0, \frac{1}{2} + 2^{-2^N}) = [0, \vartheta_{2^N-1})$ ,  $I_1 = [\frac{3}{4}, 1] = [\vartheta_1, 1]$ , and

$$I_k = [\vartheta_k, \vartheta_{k-1}) \text{ for all } 2 \leq k \leq 2^N - 1.$$

Fix  $k \in \{2, \dots, 2^N - 1\}$  and assume  $\alpha \in I_k$ . Then

$$\vartheta^{(\alpha, n)} = \arg \min_{\vartheta} \{nD(\alpha||\vartheta) + Kw(\vartheta) \ln 2\} = \arg \min_{\vartheta} \{nD(\alpha||\vartheta)\} \in \{\vartheta_k, \vartheta_{k-1}\}$$

according to (1). So clearly  $(\vartheta^{(\alpha, n)} - \vartheta_0)^2 \geq (\vartheta_k - \vartheta_0)^2 = 2^{-2k-2}$  holds. Since  $p(\alpha|n)$  decreases for increasing  $|\alpha - \vartheta_0|$ , we have  $p(\alpha|n) \geq p(\vartheta_{k-1}|n)$ . The interval  $I_k$  has length  $2^{-k-1}$ , so there are at least  $\lfloor n2^{-k-1} \rfloor \geq n2^{-k-1} - 1$  observed fractions  $\alpha$  falling in the interval. From (6), the total contribution of  $\alpha \in I_k$  can be estimated by

$$C(k) \geq \sum_{n=1}^{\infty} 2^{-2k-2}(n2^{-k-1} - 1)p(\vartheta_{k-1}|n).$$

Note that the terms in the sum even become negative for small  $n$ , which does not cause any problems. We proceed with

$$p(\vartheta_{k-1}|n) \geq \frac{1}{\sqrt{8 \cdot 2^{-2}n}} \exp[-nD(\frac{1}{2} + 2^{-k}||\frac{1}{2})] \geq \frac{1}{\sqrt{2n}} \exp[-n\frac{8}{3}2^{-2k}]$$

according to Lemma 3 and Lemma 2 (ii). By Lemma 4 (i) and (ii), we have

$$\begin{aligned} \sum_{n=1}^{\infty} \sqrt{n} \exp[-n \frac{8}{3} 2^{-2k}] &\geq \frac{\sqrt{\pi}}{2} \left(\frac{3}{8}\right)^{\frac{3}{2}} 2^{3k} - \frac{1}{\sqrt{2e}} \sqrt{\frac{3}{8}} 2^k \text{ and} \\ -\sum_{n=1}^{\infty} n^{-\frac{1}{2}} \exp[-n \frac{8}{3} 2^{-2k}] &\geq -\sqrt{\pi} \sqrt{\frac{3}{8}} 2^k. \end{aligned}$$

Considering only  $k \geq 5$ , we thus obtain

$$\begin{aligned} C(k) &\geq \frac{1}{\sqrt{2}} \sqrt{\frac{3}{8}} 2^{-2k-2} \left[ \frac{3\sqrt{\pi}}{16} 2^{2k-1} - \frac{1}{\sqrt{2e}} 2^{-1} - \sqrt{\pi} 2^k \right] \\ &\geq \frac{\sqrt{3}}{16} \left[ 3\sqrt{\pi} 2^{-5} - \frac{1}{\sqrt{2e}} 2^{-2k-1} - \sqrt{\pi} 2^{-k} \right] \geq \frac{\sqrt{3\pi}}{8} 2^{-5} - \frac{\sqrt{3}}{16\sqrt{2e}} 2^{-11} > \frac{1}{84}. \end{aligned}$$

Ignoring the contributions for  $k \leq 4$ , this implies the assertion.  $\square$

This result shows that if the parameters and their weights are chosen in an appropriate way, then the total expected error is of order  $w_0^{-1}$  instead of  $\ln w_0^{-1}$ . Interestingly, this outcome seems to depend on the arrangement and the weights of the *false* parameters rather than on the weight of the *true* one. One can check with moderate effort that the proposition still remains valid if e.g.  $w_0$  is twice as large as the other weights. Actually, the proof of Proposition 5 shows even a slightly more general result, namely admitting additional arbitrary parameters with larger complexities:

**Corollary 6** *Let  $\Theta = \{\vartheta_k : k \geq 0\}$ ,  $\vartheta_0 = \frac{1}{2}$ ,  $\vartheta_k = \frac{1}{2} + 2^{-k-1}$  for  $1 \leq k \leq 2^N - 2$ , and  $\vartheta_k \in [0, 1]$  arbitrary for  $k \geq 2^N - 1$ . Let  $Kw(\vartheta_k) = N$  for  $0 \leq k \leq 2^N - 2$  and  $Kw(\vartheta_k) > N$  for  $k \geq 2^N - 1$ . Then  $\sum_n \mathbf{E}(\vartheta_0 - \vartheta^x)^2 \geq \frac{1}{84}(2^N - 6)$  holds.*

We will use this result only for Example 16. Other and more general assertions can be proven similarly.

## 4 Upper Bounds

Although the cumulative error may be large, as seen in the previous section, the instantaneous error is always small. It is easy to demonstrate this for the Bernoulli case, to which we restrict in this paper. Much more general results have been obtained for arbitrary classes of i.i.d. models [Li99, Zha04]. Strong instantaneous bounds hold in particular if MDL is modified by replacing the factor  $\ln 2$  in (1) by something larger (e.g.  $(1 + \varepsilon) \ln 2$ ) such that complexity is penalized slightly more than usually. Note that our cumulative bounds are incomparable to these and other instantaneous bounds.

**Proposition 7** For  $n \geq 3$ , the expected instantaneous square loss is bounded as follows:

$$\mathbf{E}(\vartheta_0 - \hat{\vartheta}^{x_{1:n}})^2 \leq \frac{(\ln 2)Kw(\vartheta_0)}{2n} + \frac{\sqrt{2(\ln 2)Kw(\vartheta_0)\ln n}}{n} + \frac{6\ln n}{n}.$$

**Proof.** We give an elementary proof for the case  $\vartheta_0 \in (\frac{1}{4}, \frac{3}{4})$  only. Like in the proof of Proposition 5, we consider the contributions of different  $\alpha$  separately. By Hoeffding's inequality,  $\mathbf{P}(|\alpha - \vartheta_0| \geq \frac{c}{\sqrt{n}}) \leq 2e^{-2c^2}$  for any  $c > 0$ . Letting  $c = \sqrt{\ln n}$ , the contributions by these  $\alpha$  are thus bounded by  $\frac{2}{n^2} \leq \frac{\ln n}{n}$ .

On the other hand, for  $|\alpha - \vartheta_0| \leq \frac{c}{\sqrt{n}}$ , recall that  $\vartheta_0$  beats any  $\vartheta$  iff (3) holds. According to  $Kw(\vartheta) \geq 0$ ,  $|\alpha - \vartheta_0| \leq \frac{c}{\sqrt{n}}$ , and Lemma 2 (i) and (ii), (3) is already implied by  $|\alpha - \vartheta| \geq \sqrt{\frac{\frac{1}{2}(\ln 2)Kw(\vartheta_0) + \frac{4}{3}c^2}{n}}$ . Clearly, a contribution only occurs if  $\vartheta$  beats  $\vartheta_0$ , therefore if the opposite inequality holds. Using  $|\alpha - \vartheta_0| \leq \frac{c}{\sqrt{n}}$  again and the triangle inequality, we obtain that

$$(\vartheta - \vartheta_0)^2 \leq \frac{5c^2 + \frac{1}{2}(\ln 2)Kw(\vartheta_0) + \sqrt{2(\ln 2)Kw(\vartheta_0)c^2}}{n}$$

in this case. Since we have chosen  $c = \sqrt{\ln n}$ , this implies the assertion.  $\square$

One can improve the bound in Proposition 7 to  $\mathbf{E}(\vartheta_0 - \hat{\vartheta}^{x_{1:n}})^2 \leq \frac{Kw(\vartheta_0)}{n}$  by a refined argument, compare [BC91]. But the high-level assertion is the same: Even if the cumulative upper bound may be infinite, the instantaneous error converges rapidly to 0. Moreover, the convergence speed depends on  $Kw(\vartheta_0)$  as opposed to  $2^{Kw(\vartheta_0)}$ . Thus  $\hat{\vartheta}$  tends to  $\vartheta_0$  rapidly in probability (recall that the assertion is not strong enough to conclude almost sure convergence). The proof does not exploit  $\sum w_\vartheta \leq 1$ , but only  $w_\vartheta \leq 1$ , hence the assertion even holds for a maximum likelihood estimator (i.e.  $w_\vartheta = 1$  for all  $\vartheta \in \Theta$ ). The theorem generalizes to i.i.d. classes. For the example in Proposition 5, the instantaneous bound implies that the bulk of losses occurs very late. This does *not* hold for general (non-i.i.d.) model classes: The total loss up to time  $n$  in [PH04a, Example 9] grows linearly in  $n$ .

We will now state our main positive result that upper bounds the cumulative loss in terms of the negative logarithm of the true weight and the *arrangement* of the false parameters. The proof is similar to that of Proposition 5. We will only give the proof idea here and defer the lengthy and tedious technical details to the appendix.

Consider the cumulated sum square error  $\sum_n \mathbf{E}(\vartheta^{(\alpha,n)} - \vartheta_0)^2$ . In order to upper bound this quantity, we will partition the open unit interval  $(0, 1)$  into a sequence of intervals  $(I_k)_{k=1}^\infty$ , each of measure  $2^{-k}$ . (More precisely: Each  $I_k$  is either an interval or a union of two intervals.) Then we will estimate the contribution of each interval to the cumulated square error,

$$C(k) = \sum_{n=1}^{\infty} \sum_{\alpha \in A_n, \vartheta^{(\alpha,n)} \in I_k} p(\alpha|n)(\vartheta^{(\alpha,n)} - \vartheta_0)^2$$

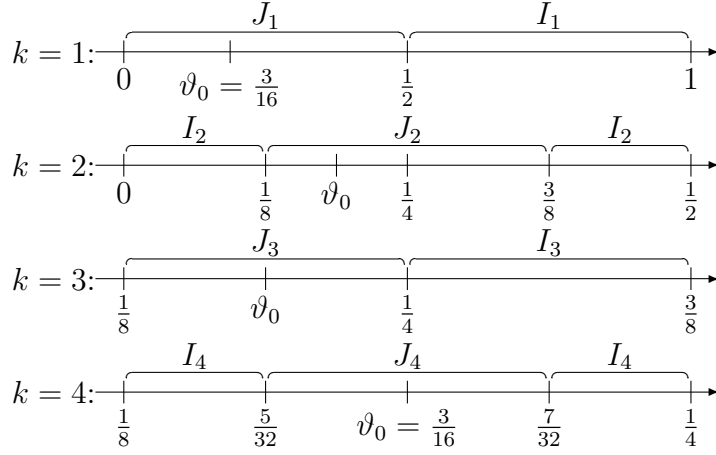


Figure 1: Example of the first four intervals for  $\vartheta_0 = \frac{3}{16}$ . We have an l-step, a c-step, an l-step and another c-step. All following steps will be also c-steps.

(compare (4) and (6)). Note that  $\vartheta^{(\alpha,n)} \in I_k$  precisely reads  $\vartheta^{(\alpha,n)} \in I_k \cap \Theta$ , but for convenience we generally assume  $\vartheta \in \Theta$  for all  $\vartheta$  being considered. This partitioning is also used for  $\alpha$ , i.e. define the contribution  $C(k, j)$  of  $\vartheta \in I_k$  where  $\alpha \in I_j$  as

$$C(k, j) = \sum_{n=1}^{\infty} \sum_{\alpha \in A_n \cap I_j, \vartheta^{(\alpha,n)} \in I_k} p(\alpha|n) (\vartheta^{(\alpha,n)} - \vartheta_0)^2.$$

We need to distinguish between  $\alpha$  that are located close to  $\vartheta_0$  and  $\alpha$  that are located far from  $\vartheta_0$ . “Close” will be roughly equivalent to  $j > k$ , “far” will be approximately  $j \leq k$ . So we get  $\sum_n \mathbf{E}(\vartheta^{(\alpha,n)} - \vartheta_0)^2 = \sum_{k=1}^{\infty} C(k) = \sum_k \sum_j C(k, j)$ . In the proof,

$$p(\alpha|n) \stackrel{\times}{\leq} [n\alpha(1-\alpha)]^{-\frac{1}{2}} \exp[-nD(\alpha||\vartheta_0)]$$

is often applied, which holds by Lemma 3 (recall that  $f \stackrel{\times}{\leq} g$  stands for  $f = O(g)$ ). Terms like  $D(\alpha||\vartheta_0)$ , arising in this context and others, can be further estimated using Lemma 2. We now give the constructions of intervals  $I_k$  and complementary intervals  $J_k$ .

**Definition 8** Let  $\vartheta_0 \in \Theta$  be given. Start with  $J_0 = [0, 1)$ . Let  $J_{k-1} = [\vartheta_k^l, \vartheta_k^r)$  and define  $d_k = \vartheta_k^r - \vartheta_k^l = 2^{-k+1}$ . Then  $I_k, J_k \subset J_{k-1}$  are constructed from  $J_{k-1}$  according to the following rules.

$$\vartheta_0 \in [\vartheta_k^l, \vartheta_k^l + \frac{3}{8}d_k) \Rightarrow J_k = [\vartheta_k^l, \vartheta_k^l + \frac{1}{2}d_k), I_k = [\vartheta_k^l + \frac{1}{2}d_k, \vartheta_k^r), \quad (7)$$

$$\vartheta_0 \in [\vartheta_k^l + \frac{3}{8}d_k, \vartheta_k^l + \frac{5}{8}d_k) \Rightarrow J_k = [\vartheta_k^l + \frac{1}{4}d_k, \vartheta_k^l + \frac{3}{4}d_k), \quad (8)$$

$$I_k = [\vartheta_k^l, \vartheta_k^l + \frac{1}{4}d_k) \cup [\vartheta_k^l + \frac{3}{4}d_k, \vartheta_k^r),$$

$$\vartheta_0 \in [\vartheta_k^l + \frac{5}{8}d_k, \vartheta_k^r) \Rightarrow J_k = [\vartheta_k^l + \frac{1}{2}d_k, \vartheta_k^r), I_k = [\vartheta_k^l, \vartheta_k^l + \frac{1}{2}d_k). \quad (9)$$

We call the  $k$ th step of the interval construction an  $l$ -step if (7) applies, a  $c$ -step if (8) applies, and an  $r$ -step if (9) applies, respectively. Fig. 1 shows an example for the interval construction.

Clearly, this is not the only possible way to define an interval construction. Maybe the reader wonders why we did not center the intervals around  $\vartheta_0$ . In fact, this construction would equally work for the proof. However, its definition would not be easier, since one still has to treat the case where  $\vartheta_0$  is located close to the boundary. Moreover, our construction has the nice property that the interval bounds are finite binary fractions.

Given the interval construction, we can identify the  $\vartheta \in I_k$  with lowest complexity:

**Definition 9** For  $\vartheta_0 \in \Theta$  and the interval construction  $(I_k, J_k)$ , let

$$\begin{aligned}\vartheta_k^I &= \arg \min \{Kw(\vartheta) : \vartheta \in I_k \cap \Theta\}, \\ \vartheta_k^J &= \arg \min \{Kw(\vartheta) : \vartheta \in J_k \cap \Theta\}, \text{ and} \\ \Delta(k) &= \max \{Kw(\vartheta_k^I) - Kw(\vartheta_k^J), 0\}.\end{aligned}$$

If there is no  $\vartheta \in I_k \cap \Theta$ , we set  $\Delta(k) = Kw(\vartheta_k^I) = \infty$ .

We can now state the main positive result of this paper. The detailed proof is deferred to the appendix. Corollaries will be given in the next section.

**Theorem 10** Let  $\Theta \subset [0, 1]$  be countable,  $\vartheta_0 \in \Theta$ , and  $w_\vartheta = 2^{-Kw(\vartheta)}$ , where  $Kw(\vartheta)$  is some complexity measure on  $\Theta$ . Let  $\Delta(k)$  be as introduced in Definition 9 and recall that  $\vartheta^x = \vartheta^{(\alpha, n)}$  depends on  $x$ 's length and observed fractions of ones. Then

$$\sum_{n=1}^{\infty} \mathbf{E}(\vartheta_0 - \vartheta^x)^2 \stackrel{\times}{\leq} Kw(\vartheta_0) + \sum_{k=1}^{\infty} 2^{-\Delta(k)} \sqrt{\Delta(k)}.$$

## 5 Uniformly Distributed Weights

We are now able to state some positive results following from Theorem 10.

**Theorem 11** Let  $\Theta \subset [0, 1]$  be a countable class of parameters and  $\vartheta_0 \in \Theta$  the true parameter. Assume that there are constants  $a \geq 1$  and  $b \geq 0$  such that

$$\min \{Kw(\vartheta) : \vartheta \in [\vartheta_0 - 2^{-k}, \vartheta_0 + 2^{-k}] \cap \Theta, \vartheta \neq \vartheta_0\} \geq \frac{k-b}{a} \quad (10)$$

holds for all  $k > aKw(\vartheta_0) + b$ . Then we have

$$\sum_{n=1}^{\infty} \mathbf{E}(\vartheta_0 - \vartheta^x)^2 \stackrel{\times}{\leq} aKw(\vartheta_0) + b \stackrel{\times}{\leq} Kw(\vartheta_0).$$

**Proof.** We have to show that

$$\sum_{k=1}^{\infty} 2^{-\Delta(k)} \sqrt{\Delta(k)} \stackrel{\times}{\leq} aKw(\vartheta_0) + b,$$

then the assertion follows from Theorem 10. Let  $k_1 = \lceil aKw(\vartheta_0) + b + 1 \rceil$  and  $k' = k - k_1$ . Then by Lemma 17 (iii) and (10) we have

$$\begin{aligned} \sum_{k=1}^{\infty} 2^{-\Delta(k)} \sqrt{\Delta(k)} &\leq \sum_{k=1}^{k_1} 1 + \sum_{k=k_1+1}^{\infty} 2^{-Kw(\vartheta_k^l) + Kw(\vartheta_0)} \sqrt{Kw(\vartheta_k^l) - Kw(\vartheta_0)} \\ &\leq k_1 + 2^{Kw(\vartheta_0)} \sum_{k=k_1+1}^{\infty} 2^{-\frac{k-b}{a}} \sqrt{\frac{k-b}{a}} \\ &\leq k_1 + 2^{Kw(\vartheta_0)} \sum_{k'=1}^{\infty} 2^{-\frac{k'+k_1-b}{a}} \sqrt{\frac{k'+k_1-b}{a}} \\ &\leq aKw(\vartheta_0) + b + 2 + \sum_{k'=1}^{\infty} 2^{-\frac{k'}{a}} \sqrt{\frac{k'}{a} + Kw(\vartheta_0)}. \end{aligned}$$

As already seen in the proof of Theorem 10,  $\sqrt{\frac{k'}{a} + Kw(\vartheta_0)} \leq \sqrt{\frac{k'}{a}} + \sqrt{Kw(\vartheta_0)}$ ,  $\sum_{k'} 2^{-\frac{k'}{a}} \stackrel{\times}{\leq} a$ , and  $\sum_{k'} 2^{-\frac{k'}{a}} \sqrt{\frac{k'}{a}} \stackrel{\times}{\leq} a$  hold. The latter is by Lemma 4 (i). This implies the assertion.  $\square$

Letting  $j = \frac{k-b}{a}$ , (10) asserts that parameters  $\vartheta$  with complexity  $Kw(\vartheta) = j$  must have a minimum distance of  $2^{-ja-b}$  from  $\vartheta_0$ . That is, if parameters with equal weights are (approximately) uniformly distributed in the neighborhood of  $\vartheta_0$ , in the sense that they are not too close to each other, then fast convergence holds. The next two results are special cases based on the set of all finite binary fractions,

$$\mathbb{Q}_{\mathbb{B}^*} = \{\vartheta = 0.\beta_1\beta_2 \dots \beta_{n-1}1 : n \in \mathbb{N}, \beta_i \in \mathbb{B}\} \cup \{0, 1\}.$$

If  $\vartheta = 0.\beta_1\beta_2 \dots \beta_{n-1}1 \in \mathbb{Q}_{\mathbb{B}^*}$ , its length is  $l(\vartheta) = n$ . Moreover, there is a binary code  $\beta'_1 \dots \beta'_{n'}$  for  $n$ , having at most  $n' \leq \lfloor \log_2(n+1) \rfloor$  bits. Then  $0\beta'_1 0\beta'_2 \dots 0\beta'_{n'} 1\beta_1 \dots \beta_{n-1}$  is a prefix-code for  $\vartheta$ . For completeness, we can define the codes for  $\vartheta = 0, 1$  to be 10 and 11, respectively. So we may define a complexity measure on  $\mathbb{Q}_{\mathbb{B}^*}$  by

$$Kw(0) = 2, Kw(1) = 2, \text{ and } Kw(\vartheta) = l(\vartheta) + 2\lfloor \log_2(l(\vartheta) + 1) \rfloor \text{ for } \vartheta \neq 0, 1. \quad (11)$$

There are other similar simple prefix codes on  $\mathbb{Q}_{\mathbb{B}^*}$  with the property  $Kw(\vartheta) \geq l(\vartheta)$ .

**Corollary 12** *Let  $\Theta = \mathbb{Q}_{\mathbb{B}^*}$ ,  $\vartheta_0 \in \Theta$  and  $Kw(\vartheta) \geq l(\vartheta)$  for all  $\vartheta \in \Theta$ , and recall  $\vartheta^x = \vartheta^{(\alpha, n)}$ . Then  $\sum_n \mathbf{E}(\vartheta_0 - \vartheta^x)^2 \stackrel{\times}{\leq} Kw(\vartheta_0)$  holds.*

**Proof.** Condition (10) holds with  $a = 1$  and  $b = 0$ .  $\square$

This is a special case of a uniform distribution of parameters with equal complexities. The next corollary is more general, it proves fast convergence if the uniform distribution is distorted by some function  $\varphi$ .

**Corollary 13** *Let  $\varphi : [0, 1] \rightarrow [0, 1]$  be an injective,  $N$  times continuously differentiable function. Let  $\Theta = \varphi(\mathbb{Q}_{\mathbb{B}^*})$ ,  $Kw(\varphi(t)) \geq l(t)$  for all  $t \in \mathbb{Q}_{\mathbb{B}^*}$ , and  $\vartheta_0 = \varphi(t_0)$  for a  $t_0 \in \mathbb{Q}_{\mathbb{B}^*}$ . Assume that there is  $n \leq N$  and  $\varepsilon > 0$  such that*

$$\begin{aligned} \left| \frac{d^n \varphi}{dt^n}(t) \right| &\geq c > 0 \quad \text{for all } t \in [t_0 - \varepsilon, t_0 + \varepsilon] \text{ and} \\ \frac{d^m \varphi}{dt^m}(t_0) &= 0 \quad \text{for all } 1 \leq m < n. \end{aligned}$$

Then we have

$$\sum \mathbf{E}(\vartheta_0 - \vartheta^x)^2 \stackrel{\times}{\leq} nKw(\vartheta_0) + 2\log_2(n!) - 2\log_2 c + n\log_2 \varepsilon \stackrel{\times}{\leq} nKw(\vartheta_0).$$

**Proof.** Fix  $j > Kw(\vartheta_0)$ , then

$$Kw(\varphi(t)) \geq j \text{ for all } t \in [t_0 - 2^{-j}, t_0 + 2^{-j}] \cap \mathbb{Q}_{\mathbb{B}^*}. \quad (12)$$

Moreover, for all  $t \in [t_0 - 2^{-j}, t_0 + 2^{-j}]$ , Taylor's theorem asserts that

$$\varphi(t) = \varphi(t_0) + \frac{d^n \varphi(\tilde{t})}{n!} (t - t_0)^n \quad (13)$$

for some  $\tilde{t}$  in  $(t_0, t)$  (or  $(t, t_0)$  if  $t < t_0$ ). We request in addition  $2^{-j} \leq \varepsilon$ , then  $|\frac{d^n \varphi}{dt^n}| \geq c$  by assumption. Apply (13) to  $t = t_0 + 2^{-j}$  and  $t = t_0 - 2^{-j}$  and define  $k = \lceil jn + \log_2(n!) - \log_2 c \rceil$  in order to obtain  $|\varphi(t_0 + 2^{-j}) - \vartheta_0| \geq 2^{-k}$  and  $|\varphi(t_0 - 2^{-j}) - \vartheta_0| \geq 2^{-k}$ . By injectivity of  $\varphi$ , we see that  $\varphi(t) \notin [\vartheta_0 - 2^{-k}, \vartheta_0 + 2^{-k}]$  if  $t \notin [t_0 - 2^{-j}, t_0 + 2^{-j}]$ . Together with (12), this implies

$$Kw(\vartheta) \geq j \geq \frac{k - \log_2(n!) + \log_2 c - 1}{n} \text{ for all } \vartheta \in [\vartheta_0 - 2^{-k}, \vartheta_0 + 2^{-k}] \cap \Theta.$$

This is condition (10) with  $a = n$  and  $b = \log_2(n!) - \log_2 c + 1$ . Finally, the assumption  $2^{-j} \leq \varepsilon$  holds if  $k \geq k_1 = n\log_2 \varepsilon + \log_2(n!) - \log_2 c + 1$ . This gives an additional contribution to the error of at most  $k_1$ .  $\square$

Corollary 13 shows an implication of Theorem 10 for *parameter identification*: A class of models is given by a set of parameters  $\mathbb{Q}_{\mathbb{B}^*}$  and a mapping  $\varphi : \mathbb{Q}_{\mathbb{B}^*} \rightarrow \Theta$ . The task is to identify the true parameter  $t_0$  or its image  $\vartheta_0 = \varphi(t_0)$ . The injectivity of  $\varphi$  is not necessary for fast convergence, but it facilitates the proof. The assumptions of Corollary 13 are satisfied if  $\varphi$  is for example a polynomial. In fact, it should be possible to prove fast convergence of MDL for many common

parameter identification problems. For sets of parameters other than  $\mathbb{Q}_{\mathbb{B}^*}$ , e.g. the set of all rational numbers  $\mathbb{Q}$ , similar corollaries can easily be proven.

How large is the constant hidden in “ $\overset{\times}{\leq}$ ”? When examining carefully the proof of Theorem 10, the resulting constant is quite huge. This is mainly due to the frequent “wasting” of small constants. The sharp bound is supposably small, perhaps 16. On the other hand, for the actual *true* expectation (as opposed to its upper bound) and complexities as in (11), numerical simulations show  $\sum_n \mathbf{E}(\vartheta_0 - \vartheta^x)^2 \leq \frac{1}{2}Kw(\vartheta_0)$ .

Finally, we state an implication which almost trivially follows from Theorem 10 but may be very useful for practical purposes, e.g. for hypothesis testing (compare [Ris99]).

**Corollary 14** *Let  $\Theta$  contain  $N$  elements,  $Kw(\cdot)$  be any complexity function on  $\Theta$ , and  $\vartheta_0 \in \Theta$ . Then we have*

$$\sum_{n=1}^{\infty} \mathbf{E}(\vartheta_0 - \vartheta^x)^2 \overset{\times}{\leq} N + Kw(\vartheta_0).$$

**Proof.**  $\sum_k 2^{-\Delta(k)} \sqrt{\Delta(k)} \leq N$  is obvious. □

## 6 The Universal Case

We briefly discuss the important universal setup, where  $Kw(\cdot)$  is (up to an additive constant) equal to the prefix Kolmogorov complexity  $K$  (that is the length of the shortest self-delimiting program printing  $\vartheta$  on some universal Turing machine). Since  $\sum_k 2^{-K(k)} \sqrt{K(k)} = \infty$  no matter how late the sum starts (otherwise there would be a shorter code for large  $k$ ), Theorem 10 does not yield a meaningful bound. This means in particular that it does not even imply our previous result, Theorem 1. But probably the following strengthening of Theorem 10 holds under the same conditions, which then easily implies Theorem 1 up to a constant.

**Conjecture 15**  $\sum_n \mathbf{E}(\vartheta_0 - \vartheta^x)^2 \overset{\times}{\leq} K(\vartheta_0) + \sum_k 2^{-\Delta(k)}$ .

Then, take an incompressible finite binary fraction  $\vartheta_0 \in \mathbb{Q}_{\mathbb{B}^*}$ , i.e.  $K(\vartheta_0) \overset{\pm}{=} l(\vartheta_0) + K(l(\vartheta_0))$ . For  $k > l(\vartheta_0)$ , we can reconstruct  $\vartheta_0$  and  $k$  from  $\vartheta_k^I$  and  $l(\vartheta_0)$  by just truncating  $\vartheta_k^I$  after  $l(\vartheta_0)$  bits. Thus  $K(\vartheta_k^I) + K(l(\vartheta_0)) \overset{\times}{\geq} K(\vartheta_0) + K(k|\vartheta_0, K(\vartheta_0))$  holds. Using Conjecture 15, we obtain

$$\sum_n \mathbf{E}(\vartheta_0 - \vartheta^x)^2 \overset{\times}{\leq} K(\vartheta_0) + 2^{K(l(\vartheta_0))} \overset{\times}{\leq} l(\vartheta_0)(\log_2 l(\vartheta_0))^2, \quad (14)$$

where the last inequality follows from the example coding given in (11).

So, under Conjecture 15, we obtain a bound which slightly exceeds the complexity  $K(\vartheta_0)$  if  $\vartheta_0$  has a certain structure. It is not obvious if the same holds for all computable  $\vartheta_0$ . In order to answer this question positive, one could try to use something like [Gác83, Eq.(2.1)]. This statement implies that as soon as  $K(k) \geq K_1$  for all  $k \geq k_1$ , we have  $\sum_{k \geq k_1} 2^{-K(k)} \stackrel{\times}{\leq} 2^{-K_1} K_1 (\log_2 K_1)^2$ . It is possible to prove an analogous result for  $\vartheta_k^I$  instead of  $k$ , however we have not found an appropriate coding that does without knowing  $\vartheta_0$ . Since the resulting bound is exponential in the code length, we therefore have not gained anything.

Another problem concerns the size of the multiplicative constant that is hidden in the upper bound. Unlike in the case of uniformly distributed weights, it is now of exponential size, i.e.  $2^{O(1)}$ . This is no artifact of the proof, as the following example shows.

**Example 16** Let  $U$  be some universal Turing machine. We construct a second universal Turing machine  $U'$  from  $U$  as follows: Let  $N \geq 1$ . If the input of  $U'$  is  $1^N p$ , where  $1^N$  is the string consisting of  $N$  ones and  $p$  is some program, then  $U$  will be executed on  $p$ . If the input of  $U'$  is  $0^N$ , then  $U'$  outputs  $\frac{1}{2}$ . Otherwise, if the input of  $U'$  is  $x$  with  $x \in \mathbb{B}^N \setminus \{0^N, 1^N\}$ , then  $U'$  outputs  $\frac{1}{2} + 2^{-x-1}$ . For  $\vartheta_0 = \frac{1}{2}$ , the conditions of Corollary 6 are satisfied (where the complexity is relative to  $U'$ ), thus  $\sum_n \mathbf{E}(\vartheta^x - \vartheta_0)^2 \stackrel{\times}{\geq} 2^N$ .

Can this also happen if the underlying universal Turing machine is not “strange” in some sense, like  $U'$ , but “natural”? Again this is not obvious. One would have to define first an appropriate notion of a “natural” universal Turing machine which rules out cases like  $U'$ . If  $N$  is of reasonable size, then one can even argue that  $U'$  is natural in the sense that its compiler constant relative to  $U$  is small.

There is a relation to the class of all *deterministic* (generally non-i.i.d.) measures. Then MDL predicts the next symbol just according to the *monotone complexity*  $Km$ , see [Hut03c]. According to [Hut03c, Theorem 5],  $2^{-Km}$  is very close to the universal semimeasure  $M$  [ZL70, Lev73]. Then the total prediction error (which is defined slightly differently in this case) can be shown to be bounded by  $2^{O(1)} Km(x_{<\infty})^3$  [Hut04]. The similarity to the (unproven) bound (14) “huge constant  $\times$  polynomial” for the universal Bernoulli case is evident.

## 7 Discussion and Conclusions

We have discovered the fact that the instantaneous and the cumulative loss bounds can be *incompatible*. On the one hand, the cumulative loss for MDL predictions may be exponential, i.e.  $2^{Kw(\vartheta_0)}$ . Thus it implies almost sure convergence at a slow speed, even for arbitrary discrete model classes [PH04a]. On the other hand, the instantaneous loss is always of order  $\frac{1}{n} Kw(\vartheta_0)$ , implying fast convergence in

probability and a cumulative loss bound of  $Kw(\vartheta_0) \ln n$ . Similar logarithmic loss bounds can be found in the literature for continuous model classes [Ris96].

A different approach to assess convergence speed is presented in [BC91]. There, an index of resolvability is introduced, which can be interpreted as the difference of the expected MDL code length and the expected code length under the true model. For discrete model classes, they show that the index of resolvability converges to zero as  $\frac{1}{n}Kw(\vartheta_0)$  [BC91, Equation (6.2)]. Moreover, they give a convergence of the predictive distributions in terms of the Hellinger distance [BC91, Theorem 4]. This implies a cumulative (Hellinger) loss bound of  $Kw(\vartheta_0) \ln n$  and therefore fast convergence in probability.

If the prior weights are arranged nicely, we have proven a small finite loss bound  $Kw(\vartheta_0)$  for MDL (Theorem 10). If parameters of equal complexity are uniformly distributed or not too strongly distorted (Theorem 11 and Corollaries), then the error is within a small multiplicative constant of the complexity  $Kw(\vartheta_0)$ . This may be applied e.g. for the case of parameter identification (Corollary 13). A similar result holds if  $\Theta$  is finite and contains only few parameters (Corollary 14), which may be e.g. satisfied for hypothesis testing. In these cases and many others, one can interpret the conditions for fast convergence as the presence of prior knowledge. One can show that if a predictor converges to the correct model, then it performs also well under arbitrarily chosen bounded loss-functions [Hut03a, Theorem 4]. From an information theoretic viewpoint one may interpret the conditions for a small bound in Theorem 10 as “good codes”.

We have proven our positive results only for Bernoulli classes, of course it would be desirable to cover more general i.i.d. classes. At least for *finite* alphabet, our assertions are likely to generalize, as this is the analog to Theorem 1 which also holds for arbitrary finite alphabet. Proving this seems even more technical than Theorem 10 and therefore not very interesting. (The interval construction has to be replaced by a sequence of nested sets in this case. Compare also the proof of the main result in [Ris96].) For small alphabets of size  $A$ , meaningful bounds can still be obtained by chaining our bounds  $A - 1$  times.

It seems more interesting to ask if our results can be *conditionalized* with respect to inputs. That is, in each time step, we are given an input and have to predict a label. This is a standard classification problem, for example a binary classification in the Bernoulli case. While it is straightforward to show that Theorem 1 still holds in this setup [PH05], it is not clear in which way the present proofs can be adapted. We leave this interesting question open.

We conclude with another open question. In abstract terms, we have proven a convergence result for the Bernoulli case by mainly exploiting the *geometry* of the space of distributions. This has been quite easy in principle, since for Bernoulli this space is just the unit interval (for i.i.d it is the space of probability vectors). It is not at all obvious if this approach can be transferred to general (computable) measures.

## A Proof of Theorem 10

The proof of Theorem 10 requires some preparations. We start by showing assertions on the interval construction from Definition 8.

**Lemma 17** *The interval construction has the following properties.*

- (i)  $|J_k| = 2^{-k}$ ,
- (ii)  $d(\vartheta_0, I_k) \geq 2^{-k-2}$ ,
- (iii)  $\max_{\vartheta \in I_k} |\vartheta - \vartheta_0| \leq 2^{-k+1}$ ,
- (iv)  $d(J_{k+5}, I_k) \geq 15 \cdot 2^{-k-6}$ .

By  $d(\cdot, \cdot)$  we mean the Euclidean distance:  $d(\tilde{\vartheta}, I) = \min\{|\tilde{\vartheta} - \vartheta| : \vartheta \in I\}$  and  $d(J, I) = \min\{d(\tilde{\vartheta}, I) : \tilde{\vartheta} \in J\}$ .

**Proof.** The first three equations are fairly obvious. The last estimate can be justified as follows. Assume that  $k$ th step of the interval construction is a c-step, the same argument applies if it is an l-step or an r-step. Let  $c$  be the center of  $J_k$  and assume without loss of generality  $\vartheta_0 \leq c$ . Define  $\vartheta_I = \max\{\vartheta \in I_k : \vartheta < c\}$  and  $\vartheta_J = \min\{\vartheta \in J_{k+5}\}$  (recall the general assumption  $\vartheta \in \Theta$  for all  $\vartheta$  that occur, i.e.  $\vartheta_I, \vartheta_J \in \Theta$ ). Then  $\vartheta_I = c - 2^{-k-1}$  and  $\vartheta_J \geq c - 2^{-k-2} - 2^{-k-6}$ , where equality holds if  $\vartheta_0 = c - 2^{-k-2}$ . Consequently,  $\vartheta_J - \vartheta_I \geq 2^{-k-1} - 2^{-k-2} - 2^{-k-6} = 15 \cdot 2^{-k-6}$ . This establishes the claim.  $\square$

Next we turn to the minimum complexity elements in the intervals.

**Proposition 18** *The following assertions hold for all  $k \geq 1$ .*

- (i)  $Kw(\vartheta_k^J) \leq Kw(\vartheta_0)$ ,
- (ii)  $Kw(\vartheta_{k+6}^J) \geq Kw(\vartheta_k^J)$ ,
- (iii)  $Kw(\vartheta_{k+1}^I) \geq Kw(\vartheta_k^J)$ ,
- (iv)  $\sum_{k=1}^{\infty} \max\{Kw(\vartheta_{k+5}^J) - Kw(\vartheta_k^I), 0\} \leq 6Kw(\vartheta_0)$ ,

**Proof.** The first three inequalities follow from  $\vartheta_0 \in J_k$  and  $J_{k+6}, I_{k+1} \subset J_k$ . This implies

$$\begin{aligned} & \sum_{j=0}^m \max\{Kw(\vartheta_{6j+6}^J) - Kw(\vartheta_{6j+1}^I), 0\} \\ & \leq \max\{Kw(\vartheta_6^J) - Kw(\vartheta_1^I), 0\} + \sum_{j=1}^m \max\{Kw(\vartheta_{6j+6}^J) - Kw(\vartheta_{6j}^J), 0\} \\ & \leq Kw(\vartheta_6^J) + \sum_{j=1}^m [Kw(\vartheta_{6j+6}^J) - Kw(\vartheta_{6j}^J)] = Kw(\vartheta_{6m+6}^J) \leq Kw(\vartheta_0) \end{aligned}$$

for all  $m \geq 0$ . By the same argument, we have

$$\sum_{j=0}^m \max \{Kw(\vartheta_{6j+i+5}^J) - Kw(\vartheta_{6j+i}^I), 0\} \leq Kw(\vartheta_0)$$

for all  $1 \leq i \leq 6$  (use (iii) in the first inequality, (ii) in the second, and (i) in the last). This implies (iv). Clearly, we could everywhere substitute 5 by some constant  $k'$  and 6 by  $k' + 1$ , but we will need the assertion only for the special case.  $\square$

Consider the case that  $\vartheta_0$  is located close to the boundary of  $[0, 1]$ . Then the interval construction involves for long time only l-steps, if we assume without loss of generality  $\vartheta_0 \leq \frac{1}{2}$ . We will need to treat this case separately, since the estimates for the general situation work only as soon as at least one c-step has taken place. Precisely, the interval construction consists only of l-steps as long as

$$\vartheta_0 < \frac{3}{4}2^{-k}, \text{ i.e. } k < -\log_2 \vartheta_0 + \log_2\left(\frac{3}{4}\right).$$

We therefore define

$$k_0 = \max \{0, \lfloor -\log_2 \vartheta_0 + \log_2 \frac{3}{4} \rfloor\} \quad (15)$$

and observe that the  $(k_0 + 1)$ st step is the first c-step. We are now prepared to give the main proof.

**Proof** of Theorem 10. Assume  $\vartheta_0 \in \Theta \setminus \{0, 1\}$ , the case  $\vartheta_0 \in \{0, 1\}$  is handled like Case 1a below and will be left to the reader.

Before we start, we will show that the contribution of  $\vartheta = 1$  to the total error is bounded by  $\frac{1}{4}$ . This is immediate, since 1 cannot become the maximizing element as soon as  $x \neq 1^n$ . Therefore the contribution is bounded by

$$\sum_{n=1}^{\infty} (1 - \vartheta_0)^2 p(1^n) = (1 - \vartheta_0)^2 \sum_{n=1}^{\infty} \vartheta_0^n = \vartheta_0(1 - \vartheta_0) \leq \frac{1}{4}. \quad (16)$$

The same is true for the contribution of  $\vartheta = 0$ .

As already mentioned, we first estimate the contributions of  $\vartheta \in I_k$  for small  $k$  if the true parameter  $\vartheta_0$  is located close to the boundary. To this aim, we assume  $\vartheta_0 \leq \frac{1}{2}$  without loss of generality. We know that the interval construction involves only l-steps as long as  $k \leq k_0$ , see (15). The very last five of these  $k$  still require a particular treatment, so we start with  $k \leq k_0 - 5$  and  $\alpha$  is far from  $\vartheta_0$ . (If  $k_0 - 5 < 1$ , then there is nothing to estimate.)

**Case 1a:**  $k \leq k_0 - 5$ ,  $j \leq k_1$ ,  $\alpha \in I_j = [2^{-j}, 2^{-j+1})$ , where  $k_1 = k + \lceil \log_2(k_0 - k - 3) \rceil + 2$ . The probability of  $\alpha$  does not exceed  $p(2^{-j})$ . The squared error may clearly be upper bounded by  $2^{-2k+2} = O(2^{-2k})$ . For  $n < 2^j$ , no such fractions can occur, so we may consider only  $n = 2^j + n'$ ,  $n' \geq 0$ . Finally, there are at most  $\lceil n \cdot 2^{-j-1} \rceil = O(2^{-j}n)$  fractions  $\alpha \in I_j$ . This follows from the general fact that if  $I \subset (0, 1)$  is any half-open or open interval of length at most  $l$ , then at most  $\lceil nl \rceil$  observed fractions can be located in  $I$ .

We now derive an estimate for the probability which is

$$p(\alpha|n) \leq p(2^{-j}|n) \stackrel{\times}{\leq} n^{-\frac{1}{2}} 2^{\frac{j}{2}} \exp[-n \cdot D(2^{-j}||\vartheta_0)]$$

according to Lemma 3. Then, Lemma 2 (v) implies

$$\exp[-nD(2^{-j}||\vartheta_0)] \leq \exp[-(2^j + n')D(2^{-j}||2^{-k_0})] \leq \exp[n'2^{-j}(k_0 - j - 1)].$$

Taking into account the upper bound for the squared error  $O(2^{-2k})$  and the maximum number of fractions  $O(2^{-j}n)$ , the contribution  $C(k, j)$  can be upper bounded by

$$C(k, j) \stackrel{\times}{\leq} \sum_{n=2^j}^{\infty} p(2^{-j}|n) 2^{-2k} \cdot 2^{-j} n \stackrel{\times}{\leq} \sum_{n'=0}^{\infty} 2^{-2k-\frac{j}{2}} \sqrt{n} \cdot \exp[n'2^{-j}(k_0 - j - 1)].$$

Decompose the right hand side using  $\sqrt{n} \leq \sqrt{2^j} + \sqrt{n'}$ . Then we have

$$\begin{aligned} \sum_{n'=0}^{\infty} 2^{-2k-\frac{j}{2}} \sqrt{2^j} \cdot \exp[n'2^{-j}(k_0 - j - 1)] &\stackrel{\times}{\leq} 2^{-2k+j}(k_0 - j - 1)^{-1} \text{ and} \\ \sum_{n'=0}^{\infty} 2^{-2k-\frac{j}{2}} \sqrt{n'} \cdot \exp[n'2^{-j}(k_0 - j - 1)] &\stackrel{\times}{\leq} 2^{-2k+j}(k_0 - j - 1)^{-\frac{3}{2}} \end{aligned}$$

where the first inequality is straightforward and the second holds by Lemma 4 (i). Letting  $k' = k_0 - k - 3$ , we have  $k' \geq 2$  and

$$(k_0 - j - 1)^{-\frac{3}{2}} \leq (k_0 - j - 1)^{-1} \leq (k_0 - k_1 - 1)^{-1} = (k' - \lceil \log_2 k' \rceil)^{-1}.$$

Thus we may conclude

$$\begin{aligned} C(k, \leq k_1) &:= \sum_{j=1}^{k_1} C(k, j) \stackrel{\times}{\leq} \sum_{j=1}^{k+\lceil \log_2 k' \rceil+2} \frac{2^{-2k+j}}{k' - \lceil \log_2 k' \rceil} \\ &\stackrel{\times}{\leq} 2^{-k} \frac{k'}{k' - \lceil \log_2 k' \rceil} \leq 2^{-k} \left( 1 + \frac{\lceil \log_2 k' \rceil}{k' - \lceil \log_2 k' \rceil} \right) \leq 3 \cdot 2^{-k} \end{aligned} \quad (17)$$

(the last inequality is sharp for  $k' = 3$ ).

**Case 1b:**  $k \leq k_0 - 5$ ,  $\alpha \leq 2^{-k_1}$  (recall  $k_1 = k + \lceil \log_2(k_0 - k - 3) \rceil + 2$ ). This means that we consider  $\alpha$  close to  $\vartheta_0$ . By (3) we know that  $\vartheta_0$  beats  $\vartheta \in I_k$  if

$$n \cdot D^\alpha(\vartheta_0||\vartheta) \geq \ln 2(Kw(\vartheta_0) - Kw(\vartheta))$$

holds. This happens certainly for  $n \geq N_1 := \ln 2 \cdot Kw(\vartheta_0) \cdot 2^{k+4}$ , since Lemma 20 below asserts  $D^\alpha(\vartheta_0||\vartheta) \geq 2^{-2-4}$ . Thus only smaller  $n$  can contribute. The total probability of all  $\alpha \leq 2^{-k_1}$  is clearly bounded by means of

$$\sum_{\alpha} p(\alpha|n) \leq 1.$$

The jump size, i.e. the squared error, is again  $O(2^{-2k})$ . Hence the total contribution caused in  $I_k$  by  $\alpha \leq 2^{-k_1}$  can thus be upper bounded by

$$C(k, > k_1) \leq \sum_{n=1}^{N_1} 2^{-2k} \leq Kw(\vartheta_0)2^{-k},$$

where  $C(k, > k_1)$  is the obvious abbreviation for this contribution. Together with (17) this implies  $C(k) \leq Kw(\vartheta_0)2^{-k}$  and therefore

$$\sum_{k=1}^{k_0-5} C(k) \leq Kw(\vartheta_0). \quad (18)$$

This finishes the estimates for  $k \leq k_0 - 5$ . We now will consider the indices

$$k_0 - 4 \leq k \leq k_0$$

and show that the contributions caused by these  $\vartheta \in I_k$  is at most  $O(Kw(\vartheta_0))$ .

**Case 2a:**  $k_0 - 4 \leq k \leq k_0$ ,  $j \leq k + 5$ ,  $\alpha \in I_j$ . Assume that  $\vartheta \in I_k$  starts contributing only for  $n > n_0$ . This is not relevant here, and we will set  $n_0 = 0$  for the moment, but then we can reuse the following computations later. Consequently we have  $n = n_0 + n'$ , and from Lemma 3 we obtain

$$p(\alpha|n) \leq n^{-\frac{1}{2}} 2^{\frac{k_0}{2}} \exp[-(n_0 + n') \cdot D(\alpha|\vartheta_0)]. \quad (19)$$

Lemma 17 implies  $d(\alpha, \vartheta_0) \geq 2^{-j-2}$  and thus

$$D(\alpha|\vartheta_0) \geq \frac{2^{-2j-4}}{2 \cdot 2^{-k_0}} = 2^{-2j-5+k_0}. \quad (20)$$

according to Lemma 2 (iii). Therefore we obtain

$$\exp[-(n_0 + n') \cdot D(\alpha|\vartheta_0)] \leq \exp[-n_0 \cdot D(\alpha|\vartheta_0)] \exp[-n' 2^{-2j-5+k_0}]. \quad (21)$$

Again the maximum square error is  $O(2^{-2k})$ , the maximum number of fractions is  $O(n2^{-j})$ . Therefore

$$C(k, j) \leq \exp[-n_0 D(\alpha|\vartheta_0)] \sum_{n'=1}^{\infty} 2^{-2k-j+\frac{k_0}{2}} \sqrt{n_0 + n'} \exp[-n' 2^{-2j-5+k_0}]. \quad (22)$$

We have

$$\sum_{n'=1}^{\infty} 2^{-2k-j+\frac{k_0}{2}} \exp[-n' 2^{-2j-5+k_0}] \leq 2^{-2k+j-\frac{k_0}{2}} \leq 2^{-2k+j} \quad \text{and} \quad (23)$$

$$\sum_{n'=1}^{\infty} 2^{-2k-j+\frac{k_0}{2}} \sqrt{n'} \exp[-n' 2^{-2j-5+k_0}] \leq 2^{-2k+2j-k_0} \leq 2^{-2k+2j}, \quad (24)$$

where the first inequality is straightforward and the second follows from Lemma 4 (i). Observe  $\sum_{j=1}^{k+5} 2^j \stackrel{\times}{\leq} 2^k$ ,  $\sum_{j=1}^{k+5} 2^{2j} \stackrel{\times}{\leq} 2^{2k}$ , and  $\sqrt{n} \leq \sqrt{n_0} + \sqrt{n'}$  in order to obtain

$$C(k, \leq k+5) \stackrel{\times}{\leq} \exp[-n_0 D(\alpha \|\vartheta_0)](1 + 2^{-k} \sqrt{n_0}). \quad (25)$$

The right hand side depends not only on  $k$  and  $n_0$ , but formally also on  $\alpha$  and even on  $\vartheta$ , since  $n_0$  itself depends on  $\alpha$  and  $\vartheta$ . Recall that for this case we have agreed on  $n_0 = 0$ , thus  $C(k, \leq k+5) = O(1)$ .

**Case 2b:**  $k_0 - 4 \leq k \leq k_0$ ,  $\alpha \in J_{k+5}$ . As before, we will argue that then  $\vartheta \in I_k$  can be the maximizing element only for small  $n$ . Namely,  $\vartheta_0$  beats  $\vartheta$  if  $n \cdot D^\alpha(\vartheta_0 \|\vartheta) \geq \ln 2(Kw(\vartheta_0) - Kw(\vartheta))$  holds. Since  $D^\alpha(\vartheta_0 \|\vartheta) \geq 2^{-2k-5}$  as stated in Lemma 20 below, this happens certainly for  $n \geq N_1 := \ln 2 \cdot Kw(\vartheta_0) \cdot 2^{2k+5}$ , thus only smaller  $n$  can contribute. Note that in order to apply Lemma 20, we need  $k \geq k_0 - 4$ . Again the total probability of all  $\alpha$  is at most 1 and the jump size is  $O(2^{-2k})$ , hence

$$C(k, > k+5) \stackrel{\times}{\leq} \sum_{n=1}^{N_1} 2^{-2k} \stackrel{\times}{\leq} Kw(\vartheta_0).$$

Together with  $C(k, \leq k+5) = O(1)$  this implies  $C(k) \stackrel{\times}{\leq} Kw(\vartheta_0)$  and thus

$$\sum_{k=k_0-4}^{k_0} C(k) \stackrel{\times}{\leq} Kw(\vartheta_0). \quad (26)$$

This completes the estimate for the initial l-steps. We now proceed with the main part of the proof. At this point, we drop the general assumption  $\vartheta_0 \leq \frac{1}{2}$ , so that we can exploit the symmetry otherwise if convenient.

**Case 3a:**  $k \geq k_0 + 1$ ,  $j \leq k + 5$ ,  $\alpha \in I_j$ . For this case, we may repeat the computations (19)-(25), arriving at

$$C(k, \leq k+5) \stackrel{\times}{\leq} \exp[-n_0 D(\alpha \|\vartheta_0)](1 + 2^{-k} \sqrt{n_0}). \quad (27)$$

The right hand side of (27) depends on  $k$  and  $n_0$  and formally also on  $\alpha$  and  $\vartheta$ . We now come to the crucial point of this proof:

*For most  $k$ ,  $n_0$  is considerably larger than 0.*

That is, for most  $k$ ,  $\vartheta \in I_k$  starts contributing late, i.e. for large  $n$ . This will cause the right hand side of (27) to be small.

We know that  $\vartheta_0$  beats  $\vartheta \in I_k$  for *any*  $\alpha \in [0, 1]$  as long as

$$nD^\alpha(\vartheta \|\vartheta_0) \leq \ln 2(Kw(\vartheta) - Kw(\vartheta_0)) \quad (28)$$

holds. We are interested in for which  $n$  this must happen regardless of  $\alpha$ , so assume that  $\alpha$  is close enough to  $\vartheta$  to make  $D^\alpha(\vartheta\|\vartheta_0) > 0$ . Since  $Kw(\vartheta) \geq Kw(\vartheta_k^I)$ , we see that (28) holds if

$$n \leq n_0(k, \alpha, \vartheta) := \frac{\ln 2 \cdot \Delta(k)}{D^\alpha(\vartheta\|\vartheta_0)}.$$

We show the following two relations:

$$\exp[-n_0(k, \alpha, \vartheta)D(\alpha\|\vartheta_0)] \leq 2^{-\Delta(k)} \quad \text{and} \quad (29)$$

$$\exp[-n_0(k, \alpha, \vartheta)D(\alpha\|\vartheta_0)]2^{-k}\sqrt{n_0(k, \alpha, \vartheta)} \stackrel{\times}{\leq} 2^{-\Delta(k)}\sqrt{\Delta(k)}, \quad (30)$$

regardless of  $\alpha$  and  $\vartheta$ . Since  $D(\alpha\|\vartheta_0) \geq D(\alpha\|\vartheta_0) - D(\alpha\|\vartheta) = D^\alpha(\vartheta\|\vartheta_0)$ , (29) is immediate. In order to verify (30), we observe that

$$D(\alpha\|\vartheta_0) \geq 2^{-2j-5+k_0} \geq 2^{-2k-15+k_0} \geq 2^{-2k-15}$$

holds as in (20). So for those  $\alpha$  and  $\vartheta$  having

$$\eta := \frac{2^{-2k-15}}{D^\alpha(\vartheta\|\vartheta_{k+5}^J)} \geq 1, \quad (31)$$

we obtain

$$\begin{aligned} \exp[-n_0(k, \alpha, \vartheta)D(\alpha\|\vartheta_0)]2^{-k}\sqrt{n_0(k, \alpha, \vartheta)} &\leq 2^{-\Delta(k)}\eta 2^{-k}\sqrt{\ln 2 \cdot \Delta(k)\eta 2^{2k+15}} \\ &\stackrel{\times}{\leq} 2^{-\Delta(k)}\sqrt{\Delta(k)}. \end{aligned}$$

since  $\eta \geq 1$ . If on the other hand (31) is not valid, then  $D^\alpha(\vartheta\|\vartheta_{k+5}^J) \stackrel{\times}{\leq} 2^{-2k}$  holds, which together with  $D(\alpha\|\vartheta_0) \geq D^\alpha(\vartheta\|\vartheta_0)$  again implies (30).

So we conclude that the dependence on  $\alpha$  and  $\vartheta$  of the right hand side of (27) is indeed only a formal one. So we obtain  $C(k, \leq k+5) \stackrel{\times}{\leq} 2^{-\Delta(k)}\sqrt{\Delta(k)}$ , hence

$$\sum_{k=k_0+1}^{\infty} C(k, \leq k+5) \stackrel{\times}{\leq} \sum_{k=1}^{\infty} 2^{-\Delta(k)}\sqrt{\Delta(k)}. \quad (32)$$

**Case 3b:**  $k \geq k_0 + 1$ ,  $\alpha \in J_{k+5}$ . We know that  $\vartheta_{k+5}^J$  beats  $\vartheta$  if

$$n \geq \ln 2 \cdot \max\{Kw(\vartheta_{k+5}^J) - Kw(\vartheta), 0\} \cdot 2^{2k+5},$$

since  $D^\alpha(\vartheta_{k+5}^J\|\vartheta) \geq 2^{-2k-5}$  according to Lemma 20. Since  $Kw(\vartheta) \geq Kw(\vartheta_k^I)$ , this happens certainly for  $n \geq N_1 := \ln 2 \cdot \max\{Kw(\vartheta_{k+5}^J) - Kw(\vartheta_k^I), 0\} \cdot 2^{2k+5}$ . Again the total probability of all  $\alpha$  is at most 1 and the jump size is  $O(2^{-2k})$ . Therefore we have

$$C(k, > k+5) \stackrel{\times}{\leq} \sum_{n=1}^{N_1} 2^{-2k} \stackrel{\times}{\leq} \max\{Kw(\vartheta_{k+5}^J) - Kw(\vartheta_k^I), 0\}.$$

Using Proposition 18 (iv), we conclude

$$\sum_{k=k_0+1}^{\infty} C(k, >k+5) \stackrel{\times}{\leq} Kw(\vartheta_0). \quad (33)$$

Combining all estimates for  $C(k)$ , namely (18), (26), (32) and (33), the assertion follows.  $\square$

**Lemma 19** *Let  $1 \leq k \leq k_0 - 5$ ,  $k_1 = k + \lceil \log_2(k_0 - k - 3) \rceil + 2$ ,  $\vartheta \geq 2^{-k}$ , and  $\alpha \leq 2^{-k_1}$ . Then  $D^\alpha(\vartheta_0 \parallel \vartheta) \geq 2^{-k-4}$  holds.*

**Proof.** By Lemma 2 (iii) and (vii), we have

$$\begin{aligned} D(\alpha \parallel \vartheta) &\geq D(2^{-k_1} \parallel 2^{-k}) \geq \frac{(2^{-k} - 2^{-k_1})^2}{2 \cdot 2^k(1 - 2^k)} \\ &\geq 2^{-k-1}(1 - 2^{-\lceil \log_2(k_0 - k - 3) \rceil - 2}) \geq 7 \cdot 2^{-k-4} \text{ and} \\ D(\alpha \parallel \vartheta_0) &\leq D(2^{-k_1} \parallel 2^{-k_0-1}) \leq 2^{-k_1}(k_0 + 1 - k_1) \\ &\leq 2^{-k-2} \frac{k_0 - k - \lceil \log_2(k_0 - k - 3) \rceil - 1}{k_0 - k - 3} \leq 6 \cdot 2^{-k-4} \end{aligned}$$

(the last inequality is sharp for  $k = k_0 - 5$ ). This implies  $D^\alpha(\vartheta_0 \parallel \vartheta) = D(\alpha \parallel \vartheta) - D(\alpha \parallel \vartheta_0) \geq 2^{-k-4}$ .  $\square$

**Lemma 20** *Let  $k \geq k_0 - 4$ ,  $\vartheta \in I_k$ , and  $\alpha, \tilde{\vartheta} \in J_{k+5}$ . Then we have  $D^\alpha(\tilde{\vartheta} \parallel \vartheta) \geq 2^{-2k-5}$ .*

**Proof.** Assume  $\vartheta \leq \frac{1}{2}$  without loss of generality. Moreover, we will only present the case  $\tilde{\vartheta} \leq \vartheta \leq \frac{1}{4}$ , the other cases are similar and simpler. From Lemma 2 (iii) and (iv) and Lemma 17 we know that

$$\begin{aligned} D(\alpha \parallel \vartheta) &\geq \frac{(\alpha - \vartheta)^2}{2\vartheta(1 - \vartheta)} \geq \frac{15^2 2^{-2k-12}}{2\vartheta} \text{ and} \\ D(\alpha \parallel \tilde{\vartheta}) &\leq \frac{3(\alpha - \tilde{\vartheta})^2}{2\alpha(1 - \alpha)} \leq \frac{4 \cdot 3 \cdot 2^{-2k-14}}{3 \cdot 2\alpha} \leq \frac{2 \cdot 128 \cdot 2^{-2k-14}}{\vartheta}. \end{aligned}$$

Note that in order to apply Lemma 2 (iv) in the second line we need to know that for  $k + 5$  a c-step has already taken place, and the last estimate follows from  $\vartheta \leq 128\alpha$  which is a consequence of  $k \geq k_0 - 4$ . Now the assertion follows from  $D^\alpha(\tilde{\vartheta} \parallel \vartheta) = D(\alpha \parallel \vartheta) - D(\alpha \parallel \tilde{\vartheta}) \geq 2^{-2k-6}(15^2 2^{-7} - 1)\vartheta^{-1} \geq 2^{-2k-5}$ .  $\square$

## References

- [BC91] A. R. Barron and T. M. Cover. Minimum complexity density estimation. *IEEE Trans. on Information Theory*, 37(4):1034–1054, 1991.
- [BRY98] A. R. Barron, J. J. Rissanen, and B. Yu. The minimum description length principle in coding and modeling. *IEEE Trans. on Information Theory*, 44(6):2743–2760, 1998.
- [CB90] B. S. Clarke and A. R. Barron. Information-theoretic asymptotics of Bayes methods. *IEEE Trans. on Information Theory*, 36:453–471, 1990.
- [Gác83] P. Gács. On the relation between descriptive complexity and algorithmic probability. *Theoretical Computer Science*, 22:71–93, 1983.
- [GL04] P. Grünwald and J. Langford. Suboptimal behaviour of Bayes and MDL in classification under misspecification. In *17th Annual Conference on Learning Theory (COLT)*, pages 331–347, 2004.
- [Hut01] M. Hutter. Convergence and error bounds for universal prediction of non-binary sequences. *Proc. 12th European Conference on Machine Learning (ECML-2001)*, pages 239–250, December 2001.
- [Hut03a] M. Hutter. Convergence and loss bounds for Bayesian sequence prediction. *IEEE Trans. on Information Theory*, 49(8):2061–2067, 2003.
- [Hut03b] M. Hutter. Optimality of universal Bayesian prediction for general loss and alphabet. *Journal of Machine Learning Research*, 4:971–1000, 2003.
- [Hut03c] M. Hutter. Sequence prediction based on monotone complexity. In *Proc. 16th Annual Conference on Learning Theory (COLT-2003)*, Lecture Notes in Artificial Intelligence, pages 506–521, Berlin, 2003. Springer.
- [Hut04] M. Hutter. Sequential predictions based on algorithmic complexity. *Journal of Computer and System Sciences*, 2005. 72(1):95–117.
- [Lev73] L. A. Levin. On the notion of a random sequence. *Soviet Math. Dokl.*, 14(5):1413–1416, 1973.
- [Li99] J. Q. Li. *Estimation of Mixture Models*. PhD thesis, Dept. of Statistics. Yale University, 1999.
- [LV97] M. Li and P. M. B. Vitányi. *An introduction to Kolmogorov complexity and its applications*. Springer, 2nd edition, 1997.
- [PH04a] J. Poland and M. Hutter. Convergence of discrete MDL for sequential prediction. In *17th Annual Conference on Learning Theory (COLT)*, pages 300–314, 2004.

- [PH04b] J. Poland and M. Hutter. On the convergence speed of MDL predictions for Bernoulli sequences. In *International Conference on Algorithmic Learning Theory (ALT)*, pages 294–308, 2004.
- [PH05] J. Poland and M. Hutter. Strong asymptotic assertions for discrete MDL in regression and classification. In *Benelearn 2005 (Ann. Machine Learning Conf. of Belgium and the Netherlands)*, 2005.
- [Ris96] J. J. Rissanen. Fisher Information and Stochastic Complexity. *IEEE Trans. on Information Theory*, 42(1):40–47, January 1996.
- [Ris99] J. J. Rissanen. Hypothesis selection and testing by the MDL principle. *The Computer Journal*, 42(4):260–269, 1999.
- [Sol78] R. J. Solomonoff. Complexity-based induction systems: comparisons and convergence theorems. *IEEE Trans. Information Theory*, IT-24:422–432, 1978.
- [VL00] P. M. Vitányi and M. Li. Minimum description length induction, Bayesianism, and Kolmogorov complexity. *IEEE Trans. on Information Theory*, 46(2):446–464, 2000.
- [Vov97] V. G. Vovk. Learning about the parameter of the Bernoulli model. *Journal of Computer and System Sciences*, 55:96–104, 1997.
- [Zha04] T. Zhang. On the convergence of MDL density estimation. In *Proc. 17th Annual Conference on Learning Theory (COLT)*, pages 315–330, 2004.
- [ZL70] A. K. Zvonkin and L. A. Levin. The complexity of finite objects and the development of the concepts of information and randomness by means of the theory of algorithms. *Russian Mathematical Surveys*, 25(6):83–124, 1970.