

---

# Monotone Conditional Complexity Bounds on Future Prediction Errors

---

Alexey Chernov and Marcus Hutter

IDSIA, Galleria 2, CH-6928 Manno-Lugano, Switzerland\*  
{alexey,marcus}@idsia.ch, <http://www.idsia.ch/~{alexey,marcus}>

18 July 2005

## Abstract

We bound the future loss when predicting any (computably) stochastic sequence online. Solomonoff finitely bounded the total deviation of his universal predictor  $M$  from the true distribution  $\mu$  by the algorithmic complexity of  $\mu$ . Here we assume we are at a time  $t > 1$  and already observed  $x = x_1 \dots x_t$ . We bound the future prediction performance on  $x_{t+1}x_{t+2}\dots$  by a new variant of algorithmic complexity of  $\mu$  given  $x$ , plus the complexity of the randomness deficiency of  $x$ . The new complexity is monotone in its condition in the sense that this complexity can only decrease if the condition is prolonged. We also briefly discuss potential generalizations to Bayesian model classes and to classification problems.

## Keywords

Kolmogorov complexity, posterior bounds, online sequential prediction, Solomonoff prior, monotone conditional complexity, total error, future loss, randomness deficiency.

---

\*This work was supported by SNF grants 200020-100259 (to Jürgen Schmidhuber), 2100-67712 and 200020-107616.

# 1 Introduction

We consider the problem of online=sequential predictions. We assume that the sequences  $x = x_1x_2x_3\dots$  are drawn from some “true” but unknown probability distribution  $\mu$ . Bayesians proceed by considering a class  $\mathcal{M}$  of models=hypotheses=distributions, sufficiently large such that  $\mu \in \mathcal{M}$ , and a prior over  $\mathcal{M}$ . Solomonoff considered the truly large class that contains all computable probability distributions [Sol64]. He showed that his universal distribution  $M$  converges rapidly to  $\mu$  [Sol78], i.e. predicts well in any environment as long as it is computable or can be modeled by a computable probability distribution (all physical theories are of this sort).  $M(x)$  is roughly  $2^{-K(x)}$ , where  $K(x)$  is the length of the shortest description of  $x$ , called Kolmogorov complexity of  $x$ . Since  $K$  and  $M$  are incomputable, they have to be approximated in practice. See e.g. [Sch02b, Hut04, LV97, CV05] and references therein. The universality of  $M$  also precludes useful statements of the prediction quality at particular time instances  $n$  [Hut04, p62], as opposed to simple classes like i.i.d. sequences (data) of size  $n$ , where accuracy is typically  $O(n^{-1/2})$ . Luckily, bounds on the expected *total*=cumulative loss (e.g. number of prediction errors) for  $M$  can be derived [Sol78, Hut03a, Hut03b], which is often sufficient in an online setting. The bounds are in terms of the (Kolmogorov) complexity of  $\mu$ . For instance, for deterministic  $\mu$ , the number of errors is (in a sense tightly) bounded by  $K(\mu)$  which measures in this case the information (in bits) in the observed infinite sequence  $x$ .

**What’s new.** In this paper we assume we are at a time  $t > 1$  and already observed  $x = x_1\dots x_t$ . Hence we are interested in the future prediction performance on  $x_{t+1}x_{t+2}\dots$ , since typically we don’t care about past errors. If the total loss is finite, the future loss must necessarily be small for large  $t$ . In a sense the paper intends to quantify this apparent triviality. If the complexity of  $\mu$  bounds the total loss, a natural guess is that something like the conditional complexity of  $\mu$  given  $x$  bounds the future loss. (If  $x$  contains a lot of (or even all) information about  $\mu$ , we should make fewer (no) errors anymore.) Indeed, we prove two bounds of this kind but with additional terms describing structural properties of  $x$ . These additional terms appear since the total loss is bounded only in expectation, and hence the future loss is small only for “most”  $x_1\dots x_t$ . In the first bound (Theorem 1), the additional term is the complexity of the length of  $x$  (a kind of worst-case estimation). The second bound (Theorem 7) is finer: the additional term is the complexity of the randomness deficiency of  $x$ . The advantage is that the deficiency is small for “typical”  $x$  and bounded on average (in contrast to the length). But in this case the conventional conditional complexity turned out to be unsuitable. So we introduce a new natural modification of conditional Kolmogorov complexity, which is monotone as a function of condition. Informally speaking, we require programs (=descriptions) to be consistent in the sense that if a program generates some  $\mu$  given  $x$ , then it must generate the same  $\mu$  given any prolongation of  $x$ . The new posterior bounds also significantly improve the previous total bounds.

**Contents.** The paper is organized as follows. Some basic notation and definitions

are given in Sections 2 and 3. In Section 4 we prove and discuss the length-based bound Theorem 1. In Section 5 we show why a new definition of complexity is necessary and formulate the deficiency-based bound Theorem 7. We discuss the definition and basic properties of the new complexity in Section 6, and prove Theorem 7 in Section 7. We briefly discuss potential generalizations to general model classes  $\mathcal{M}$  and classification in the concluding Section 8.

## 2 Notation & Definitions

We essentially follow the notation of [LV97, Hut04].

**Strings and natural numbers.** We write  $\mathcal{X}^*$  for the set of finite strings over a finite alphabet  $\mathcal{X}$ , and  $\mathcal{X}^\infty$  for the set of infinite sequences. The cardinality of a set  $\mathcal{S}$  is denoted by  $|\mathcal{S}|$ . We use letters  $i, k, l, n, t$  for natural numbers,  $u, v, x, y, z$  for finite strings,  $\epsilon$  for the empty string, and  $\alpha = \alpha_{1:\infty}$  etc. for infinite sequences. For a string  $x$  of length  $\ell(x) = n$  we write  $x_1x_2\dots x_n$  with  $x_t \in \mathcal{X}$  and further abbreviate  $x_{k:n} := x_kx_{k+1}\dots x_{n-1}x_n$  and  $x_{<n} := x_1\dots x_{n-1}$ . For  $x_t \in \mathcal{X}$ , denote by  $\bar{x}_t$  an arbitrary element from  $\mathcal{X}$  such that  $\bar{x}_t \neq x_t$ . For binary alphabet  $\mathcal{X} = \{0,1\}$ , the  $\bar{x}_t$  is uniquely defined. We occasionally identify strings with natural numbers.

**Prefix sets.** A string  $x$  is called a (proper) prefix of  $y$  if there is a  $z (\neq \epsilon)$  such that  $xz = y$ ;  $y$  is called a prolongation of  $x$ . We write  $x* = y$  in this case, where  $*$  is a wildcard for a string, and similarly for infinite sequences. A set of strings is called prefix free if no element is a proper prefix of another. Any prefix set  $\mathcal{P}$  has the important property of satisfying Kraft's inequality  $\sum_{x \in \mathcal{P}} |\mathcal{X}|^{-\ell(x)} \leq 1$ .

**Asymptotic notation.** We write  $f(x) \stackrel{\times}{\sim} g(x)$  for  $f(x) = O(g(x))$  and  $f(x) \stackrel{\pm}{\sim} g(x)$  for  $f(x) \leq g(x) + O(1)$ . Equalities  $\stackrel{\times}{\sim}$ ,  $\stackrel{\pm}{\sim}$  are defined similarly: they hold if the corresponding inequalities hold in both directions.

**(Semi)measures.** We call  $\rho: \mathcal{X}^* \rightarrow [0,1]$  a (semi)measure iff  $\sum_{x_n \in \mathcal{X}} \rho(x_{1:n}) \stackrel{(\leq)}{=} \rho(x_{<n})$  and  $\rho(\epsilon) \stackrel{(\leq)}{=} 1$ .  $\rho(x)$  is interpreted as the  $\rho$ -probability of sampling a sequence which starts with  $x$ . The conditional probability (posterior)  $\rho(y|x) := \frac{\rho(xy)}{\rho(x)}$  is the  $\rho$ -probability that a string  $x$  is followed by (continued with)  $y$ . We call  $\rho$  deterministic if  $\exists \alpha: \rho(\alpha_{1:n}) = 1 \forall n$ . In this case we identify  $\rho$  with  $\alpha$ .

**Random events and expectations.** We assume that sequence  $\omega = \omega_{1:\infty}$  is sampled from the "true" measure  $\mu$ , i.e.  $\mathbf{P}[\omega_{1:n} = x_{1:n}] = \mu(x_{1:n})$ . We denote expectations w.r.t.  $\mu$  by  $\mathbf{E}$ , i.e. for a function  $f: \mathcal{X}^n \rightarrow \mathbb{R}$ ,  $\mathbf{E}[f] = \mathbf{E}[f(\omega_{1:n})] = \sum_{x_{1:n}} \mu(x_{1:n}) f(x_{1:n})$ . We abbreviate  $\mu_t := \mu(x_t | \omega_{<t})$ .

**Enumerable sets and functions.** A set of strings (or naturals, or other constructive objects) is called *enumerable* if it is the range of some computable function. A function  $f: \mathcal{X}^* \rightarrow \mathbb{R}$  is called *(co-)enumerable* if the set of pairs  $\{\langle x, \frac{k}{n} \rangle \mid f(x) \stackrel{(\leq)}{>} \frac{k}{n}\}$  is enumerable. A measure  $\mu$  is called *computable* if it is enumerable and co-enumerable and the set  $\{x \mid \mu(x) = 0\}$  is decidable (i.e. enumerable and co-enumerable).

**Prefix Kolmogorov complexity.** The conditional prefix complexity  $K(y|x) := \min\{\ell(p) : U(p,x) = y\}$  is the length of the shortest binary (self-delimiting) program  $p \in \{0,1\}^*$  on a universal prefix Turing machine  $U$  with output  $y \in \mathcal{X}^*$  and input  $x \in \mathcal{X}^*$  [LV97].  $K(x) := K(x|\epsilon)$ . For non-string objects  $o$  we define  $K(o) := K(\langle o \rangle)$ , where  $\langle o \rangle \in \mathcal{X}^*$  is some standard code for  $o$ . In particular, if  $(f_i)_{i=1}^\infty$  is an enumeration of all (co-)enumerable functions, we define  $K(f_i) := K(i)$ . We need the following properties: The co-enumerability of  $K$ , the upper bounds  $K(x|\ell(x)) \stackrel{\pm}{\leq} \ell(x)\log_2|\mathcal{X}|$  and  $K(n) \stackrel{\pm}{\leq} 2\log_2 n$ , Kraft's inequality  $\sum_x 2^{-K(x)} \leq 1$ , the lower bound  $K(x) \geq l(x)$  for “most”  $x$  (which implies  $K(n) \xrightarrow{n \rightarrow \infty} \infty$ ), extra information bounds  $K(x|y) \stackrel{\pm}{\leq} K(x) \stackrel{\pm}{\leq} K(x,y)$ , subadditivity  $K(xy) \stackrel{\pm}{\leq} K(x,y) \stackrel{\pm}{\leq} K(y) + K(x|y)$ , information non-increase  $K(f(x)) \stackrel{\pm}{\leq} K(x) + K(f)$  for computable  $f: \mathcal{X}^* \rightarrow \mathcal{X}^*$ , and coding relative to a probability distribution (MDL): if  $P: \mathcal{X}^* \rightarrow [0,1]$  is enumerable and  $\sum_x P(x) \leq 1$ , then  $K(x) \stackrel{\pm}{\leq} -\log_2 P(x) + K(P)$ .

**Monotone and Solomonoff complexity.** The monotone complexity  $Km(x) := \min\{\ell(p) : U(p) = x^*\}$  is the length of the shortest binary (possibly non-halting) program  $p \in \{0,1\}^*$  on a universal monotone Turing machine  $U$  which outputs a string starting with  $x$ . Solomonoff's prior  $M(x) := \sum_{p:U(p)=x^*} 2^{-\ell(p)} =: 2^{-KM(x)}$  is the probability that  $U$  outputs a string starting with  $x$  if provided with fair coin flips on the input tape. Most complexities coincide within an additive term  $O(\log \ell(x))$ , e.g.  $K(x|\ell(x)) \stackrel{\pm}{\leq} KM(x) \leq Km(x) \leq K(x)$ , hence similar relations as for  $K$  hold.

### 3 Setup

**Convergent predictors.** We assume that  $\mu$  is a “true”<sup>1</sup> sequence generating measure, also called environment. If we know the generating process  $\mu$ , and given past data  $x_{<t}$ , we can predict the probability  $\mu(x_t|x_{<t})$  of the next data item  $x_t$ . Usually we do not know  $\mu$ , but estimate it from  $x_{<t}$ . Let  $\rho(x_t|x_{<t})$  be an estimated probability<sup>2</sup> of  $x_t$ , given  $x_{<t}$ . Closeness of  $\rho(x_t|x_{<t})$  to  $\mu(x_t|x_{<t})$  is desirable as a goal in itself or when performing a Bayes decision  $y_t$  that has minimal  $\rho$ -expected loss  $l_t^\rho(x_{<t}) := \min_{y_t} \sum_{x_t} \text{Loss}(x_t, y_t) \rho(x_t|x_{<t})$ . Consider, for instance, a weather data sequence  $x_{1:n}$  with  $x_t = 1$  meaning rain and  $x_t = 0$  meaning sun at day  $t$ . Given  $x_{<t}$  the probability of rain tomorrow is  $\mu(1|x_{<t})$ . A weather forecaster may announce the probability of rain to be  $y_t := \rho(1|x_{<t})$ , which should be close to the true probability  $\mu(1|x_{<t})$ . To aim for

$$\rho(x'_t|x_{<t}) - \mu(x'_t|x_{<t}) \xrightarrow{(fast)} 0 \quad \text{for } t \rightarrow \infty$$

seems reasonable.

**Convergence in mean sum.** We can quantify the deviation of  $\rho_t$  from  $\mu_t$ , e.g. by the squared difference

$$s_t(\omega_{<t}) := \sum_{x_t \in \mathcal{X}} (\rho(x_t|\omega_{<t}) - \mu(x_t|\omega_{<t}))^2 \equiv \sum_{x_t} (\rho_t - \mu_t)^2$$

<sup>1</sup>Also called *objective* or *aleatory* probability or *chance*.

<sup>2</sup>Also called *subjective* or *belief* or *epistemic* probability.

Alternatively one may also use the squared absolute distance  $s_t := \frac{1}{2}(\sum_{x_t} |\rho_t - \mu_t|)^2$ , the Hellinger distance  $s_t := \sum_{x_t} (\sqrt{\rho_t} - \sqrt{\mu_t})^2$ , the KL-divergence  $s_t := \sum_{x_t} \mu_t \ln \frac{\mu_t}{\rho_t}$ , or the squared Bayes regret  $s_t := \frac{1}{2}(l_t^\rho - l_t^\mu)^2$  for  $l_t \in [0,1]$ . For all these distances one can show [Hut03a, Hut04] that their cumulative expectation from  $l$  to  $n$  is bounded as follows:

$$0 \leq \mathbf{E}\left[\sum_{t=l}^n s_t | \omega_{<l}\right] \leq \mathbf{E}\left[\ln \frac{\mu(\omega_{l:n} | \omega_{<l})}{\rho(\omega_{l:n} | \omega_{<l})} | \omega_{<l}\right] =: D_{l:n}(\omega_{<l}). \quad (1)$$

$D_{l:n}$  is increasing in  $n$ , hence  $D_{l:\infty} \in [0, \infty]$  exists [Hut01, Hut04]. A sequence of random variables like  $s_t$  is said to converge to zero with probability 1 if the set  $\{\omega : s_t(\omega) \xrightarrow{t \rightarrow \infty} 0\}$  has measure 1.  $s_t$  is said to converge to zero in mean sum if  $\sum_{t=1}^{\infty} \mathbf{E}[|s_t|] \leq c < \infty$ , which implies convergence with probability 1 (rapid if  $c$  is of reasonable size). Therefore a small finite bound on  $D_{1:\infty}$  would imply rapid convergence of the  $s_t$  defined above to zero, hence  $\rho_t \rightarrow \mu_t$  and  $l_t^\rho \rightarrow l_t^\mu$  fast. So the crucial quantities to consider and bound (in expectation) are  $\ln \frac{\mu(x)}{\rho(x)}$  if  $l=1$  and  $\ln \frac{\mu(y|x)}{\rho(y|x)}$  for  $l>1$ . For illustration we will sometimes loosely interpret  $D_{1:\infty}$  and other quantities as the number of prediction errors, as for the error-loss they are closely related to it [Hut01].

**Bayes mixtures.** A Bayesian considers a class of distributions  $\mathcal{M} := \{\nu_1, \nu_2, \dots\}$ , large enough to contain  $\mu$ , and uses the Bayes mixture

$$\xi(x) := \sum_{\nu \in \mathcal{M}} w_\nu \cdot \nu(x), \quad \sum_{\nu \in \mathcal{M}} w_\nu = 1, \quad w_\nu > 0. \quad (2)$$

for prediction, where  $w_\nu$  can be interpreted as the prior of (or initial belief in)  $\nu$ . The dominance

$$\xi(x) \geq w_\mu \cdot \mu(x) \quad \forall x \in \mathcal{X}^* \quad (3)$$

is its most important property. Using  $\rho = \xi$  for prediction, this implies  $D_{1:\infty} \leq \ln w_\mu^{-1} < \infty$ , hence  $\xi_t \rightarrow \mu_t$ . If  $\mathcal{M}$  is chosen sufficiently large, then  $\mu \in \mathcal{M}$  is not a serious constraint.

**Solomonoff prior.** So we consider the largest (from a computational point of view) relevant class, the class  $\mathcal{M}_U$  of all enumerable semimeasures (which includes all computable probability distributions) and choose  $w_\nu = 2^{-K(\nu)}$  which is biased towards simple environments (Occam's razor). This gives us Solomonoff-Levin's prior  $M$  [Sol64, ZL70] (this definition coincides within an irrelevant multiplicative constant with the one in Section 2). In the following we assume  $\mathcal{M} = \mathcal{M}_U$ ,  $\rho = \xi = M$ ,  $w_\nu = 2^{-K(\nu)}$  and  $\mu \in \mathcal{M}_U$  being a computable (proper) measure, hence  $M(x) \geq 2^{-K(\mu)} \mu(x) \forall x$  by (3).

**Prediction of deterministic environments.** Consider a computable sequence  $\alpha = \alpha_{1:\infty}$  "sampled from  $\mu \in \mathcal{M}$ " with  $\mu(\alpha) = 1$ , i.e.  $\mu$  is deterministic, then from (3) we get

$$\sum_{t=1}^{\infty} |1 - M(\alpha_t | \alpha_{<t})| \leq -\sum_{t=1}^{\infty} \ln M(\alpha_t | \alpha_{<t}) = -\ln M(\alpha_{1:\infty}) \leq K(\mu) \ln 2 < \infty, \quad (4)$$

which implies that  $M(\alpha_t|\alpha_{<t})$  converges rapidly to 1 and hence  $M(\bar{\alpha}_t|\alpha_{<t}) \rightarrow 0$ , i.e. asymptotically  $M$  correctly predicts the next symbol. The number of prediction errors is of the order of the complexity  $K(\mu) \stackrel{\pm}{=} Km(\alpha)$  of the sequence.

For binary alphabet this is the best we can expect, since at each time-step only a single bit can be learned about the environment, and only after we “know” the environment we can predict correctly. For non-binary alphabet,  $K(\mu)$  still measures the information in  $\mu$  in bits, but feedback per step can now be  $\log_2|\mathcal{X}|$  bits, so we may expect a better bound  $K(\mu)/\log_2|\mathcal{X}|$ . But in the worst case all  $\alpha_t \in \{0,1\} \subseteq \mathcal{X}$ . So without structural assumptions on  $\mu$  the bound cannot be improved even if  $\mathcal{X}$  is huge. We will see how our posterior bounds can help in this situation.

**Individual randomness (deficiency).** Let us now consider a general (not necessarily deterministic) computable measure  $\mu \in \mathcal{M}$ . The Shannon-Fano code of  $x$  w.r.t.  $\mu$  has code-length  $\lceil -\log_2\mu(x) \rceil$ , which is “optimal” for “typical/random”  $x$  sampled from  $\mu$ . Further,  $-\log_2M(x) \approx K(x)$  is the length of an “optimal” code for  $x$ . Hence  $-\log_2\mu(x) \approx -\log_2M(x)$  for “ $\mu$ -typical/random”  $x$ . This motivates the definition of  $\mu$ -randomness deficiency

$$d_\mu(x) := \log_2 \frac{M(x)}{\mu(x)}$$

which is small for “typical/random”  $x$ . Formally, a sequence  $\alpha$  is called (Martin-Löf) random iff  $d_\mu(\alpha) := \sup_n d_\mu(\alpha_{1:n}) < \infty$ , i.e. iff its Shannon-Fano code is “optimal” (note that  $d_\mu(\alpha) \geq -K(\mu) > -\infty$  for all sequences), i.e. iff

$$\sup_n \left| \sum_{t=1}^n \log \frac{\mu(\alpha_t|\alpha_{<t})}{M(\alpha_t|\alpha_{<t})} \right| \equiv \sup_n \left| \log \frac{\mu(\alpha_{1:n})}{M(\alpha_{1:n})} \right| < \infty.$$

Unfortunately this does not imply  $M_t \rightarrow \mu_t$  on the  $\mu$ -random  $\alpha$ , since  $M_t$  may oscillate around  $\mu_t$ , which indeed can happen [HM04]. But if we take the expectation, Solomonoff [Sol78, Hut01, Hut04] showed

$$0 \leq \sum_{t=1}^{\infty} \mathbf{E} \sum_{x_t} (M_t - \mu_t)^2 \leq D_{1:\infty} = \lim_{n \rightarrow \infty} \mathbf{E}[-d_\mu(\omega_{1:n})] \ln 2 \leq K(\mu) \ln 2 < \infty \quad (5)$$

hence,  $M_t \rightarrow \mu_t$  with  $\mu$ -probability 1. So in any case,  $d_\mu(x)$  is an important quantity, since the smaller  $-d_\mu(x)$  (at least in expectation) the better  $M$  predicts.

## 4 Posterior Bounds

**Posterior bounds.** Both bounds, (4) and (5) bound the total (cumulative) discrepancy (error) between  $M_t$  and  $\mu_t$ . Since the discrepancy sum  $D_{1:\infty}$  is finite, we know that after sufficiently long time  $t=l$ , we will make little further errors, i.e. the future error sum  $D_{l:\infty}$  is small. The main goal of this paper is to quantify this asymptotic statement. So we need bounds on  $\log_2 \frac{\mu(y|x)}{M(y|x)}$ , where  $x$  are past and  $y$  are future

observations. Since  $\log_2 \frac{\mu(y)}{M(y)} \leq K(\mu)$  and  $\mu(y|x)/M(y|x)$  are conditional versions of true/universal distributions, it seems natural that the unconditional bound  $K(\mu)$  also simply conditionalizes to  $\log_2 \frac{\mu(y|x)}{M(y|x)} \stackrel{?}{\leq} K(\mu|x)$ . The more information the past observation  $x$  contains about  $\mu$ , the easier it is to code  $\mu$  i.e. the smaller is  $K(\mu|x)$ , and hence the less future predictions errors  $D_{l:\infty}$  we should make. Once  $x$  contains all information about  $\mu$ , i.e.  $K(\mu|x) \stackrel{\pm}{=} 0$ , we should make no errors anymore. More formally, optimally coding  $x$  then  $\mu|x$  and finally  $y|\mu, x$  by Shannon-Fano, gives a code for  $xy$ , hence  $K(xy) \lesssim K(x) + K(\mu|x) + \log_2 \mu(y|x)^{-1}$ . Since  $K(z) \approx -\log_2 M(z)$  this implies  $\log_2 \frac{\mu(y|x)}{M(y|x)} \lesssim K(\mu|x)$ , but with logarithmic fudge that tends to infinity for  $\ell(y) \rightarrow \infty$ , which is unacceptable. The  $y$ -independent bound we need was first stated in [Hut04, Prob.2.6(iii)]:

**Theorem 1.** *For any computable measure  $\mu$  and any  $x, y \in \mathcal{X}^*$  it holds*

$$\log_2 \frac{\mu(y|x)}{M(y|x)} \stackrel{\pm}{\leq} K(\mu|x) + K(\ell(x)).$$

**Proof.** For any fixed  $l$  we define the following function of  $z \in \mathcal{X}^*$ . For  $\ell(z) \geq l$ ,

$$\psi_l(z) := \sum_{\nu \in \mathcal{M}} 2^{-K(\nu|z_{1:l})} M(z_{1:l}) \nu(z_{l+1:\ell(z)}).$$

For  $\ell(z) < l$  we extend  $\psi_l$  by defining  $\psi_l(z) := \sum_{u: \ell(u) = l - \ell(z)} \psi_l(zu)$ . It is easy to see that  $\psi_l$  is an enumerable semimeasure. By definition of  $M$ , we have  $M(z) \geq 2^{-K(\psi_l)} \psi_l(z)$  for any  $l$  and  $z$ . Now let  $l = \ell(x)$  and  $z = xy$ . Let us define a semimeasure  $\mu_x(y) := \mu(y|x)$ . Then

$$M(xy) \geq 2^{-K(\psi_l)} \psi_l(xy) \geq 2^{-K(\psi_l)} 2^{-K(\mu_x|x)} M(x) \mu_x(y).$$

Taking the logarithm, after trivial transformations, we get  $\log_2 \frac{\mu(y|x)}{M(y|x)} \leq K(\mu_x|x) + K(\psi_l)$ . To complete the proof, let us note that  $K(\psi_l) \stackrel{\pm}{\leq} K(l)$  and  $K(\mu_x|x) \stackrel{\pm}{\leq} K(\mu|x)$ .  $\square$

**Corollary 2.** *The future and total deviations of  $M_t$  from  $\mu_t$  are bounded by*

$$\begin{aligned} i) \quad & \sum_{t=l+1}^{\infty} \mathbf{E}[s_t | \omega_{1:l}] \leq D_{l+1:\infty}(\omega_{1:l}) \stackrel{\pm}{\leq} (K(\mu|\omega_{1:l}) + K(l)) \ln 2 \\ ii) \quad & \sum_{t=1}^{\infty} \mathbf{E}[s_t] \stackrel{\pm}{\leq} \min_l \{ \mathbf{E}[K(\mu|\omega_{1:l}) + K(l)] \ln 2 + 2l \} \end{aligned}$$

**Proof.** (i) The first inequality is (1) and the second follows by taking the conditional expectation  $\mathbf{E}[\cdot | \omega_{1:l}]$  in Theorem 1. (ii) follows from (i) by taking the unconditional expectation and from  $\sum_{t=1}^l \mathbf{E}[s_t] \leq 2l$ , since  $s_t \leq 2$ .  $\square$

**Examples and more motivation.** The bounds Theorem 1 and Corollary 2(i) prove and quantify the intuition that the more we know about the environment, the better our predictions. We show the usefulness of the new bounds for some deterministic environments  $\mu \hat{=} \alpha$ .

Assume all observations are identical, i.e.  $\alpha = x_1 x_1 x_1 \dots$ . Further assume that  $\mathcal{X}$  is huge and  $K(x_1) = \log_2 |\mathcal{X}|$ , i.e.  $x_1$  is a typical/random/complex element of  $\mathcal{X}$ . For instance if  $x_1$  is a  $256^3$  color  $512 \times 512$  pixel image, then  $|\mathcal{X}| = 256^{3 \times 512 \times 512}$ . Hence the standard bound (5) on the number of errors  $D_{1:\infty} / \ln 2 \leq K(\mu) \stackrel{\pm}{=} K(x_1) = 3 \cdot 2^{21}$  is huge. Of course, interesting pictures are not purely random, but their complexity is often only a factor 10..100 less, so still large. On the other hand, any reasonable prediction scheme observing a few (rather than several thousands) identical images, should predict that the next image will be the same. This is what our posterior bound gives,  $D_{2:\infty}(x_1) \stackrel{\pm}{=} K(\mu|x_1) + K(1) \stackrel{\pm}{=} 0$ , hence indeed  $M$  makes only  $\sum_{t=1}^{\infty} \mathbf{E}[s_t] = O(1)$  errors by Corollary 2(ii), significantly improving upon Solomonoff's bound  $K(\mu) \ln 2$ .

More generally, assume  $\alpha = x\omega$ , where the initial part  $x = x_{1:l}$  contains all information about the remainder, i.e.  $K(\mu|x) \stackrel{\pm}{=} K(\omega|x) \stackrel{\pm}{=} 0$ . For instance,  $x$  may be a binary program for  $\pi$  or  $e$  and  $\omega$  be its  $|\mathcal{X}|$ -ary expansion. Sure, given the algorithm for some number sequence, it should be perfectly predictable. Indeed, Theorem 1 implies  $D_{l+1:\infty} \stackrel{\pm}{=} K(l)$ , which can be exponentially smaller than Solomonoff's bound  $K(\mu)$  ( $\stackrel{\pm}{=} l$  if  $K(x) \stackrel{\pm}{=} \ell(x)$ ). On the other hand,  $K(l) \geq \log_2 l$  for most  $l$ , i.e. is larger than  $O(1)$  what one might hope for.

**Logarithmic versus constant accuracy.** So there is one blemish in the bound. There is an additive correction of logarithmic size in the length of  $x$ . Many theorems in algorithmic information theory hold to within an additive constant, sometimes this is easily reached, sometimes hard, sometimes one needs a suitable complexity variant, and sometimes the logarithmic accuracy cannot be improved [LV97]. The latter is the case with Theorem 1:

**Lemma 3.** *For  $\mathcal{X} = \{0,1\}$ , for any computable measure  $\mu$ , there exists a computable sequence  $\alpha \in \{0,1\}^{\infty}$  such that for any  $l \in \mathbb{N}$*

$$D_{l:\infty}(\alpha_{<l}) \geq D_{l:l}(\alpha_{<l}) \equiv \sum_{b \in \{0,1\}} \mu(b|\alpha_{<l}) \ln \frac{\mu(b|\alpha_{<l})}{M(b|\alpha_{<l})} \stackrel{\pm}{=} \frac{1}{3} K(l).$$

**Proof.** Let us construct a computable sequence  $\alpha \in \{0,1\}^{\infty}$  by induction. Assume that  $\alpha_{<l}$  is constructed. Since  $\mu$  is a measure, either  $\mu(0|\alpha_{<l}) > c$  or  $\mu(1|\alpha_{<l}) > c$  for  $c := [3 \ln 2]^{-1} < \frac{1}{2}$ . Since  $\mu$  is computable, we can find (effectively)  $b \in \{0,1\}$  such that  $\mu(b|\alpha_{<l}) > c$ . Put  $\alpha_l = \bar{b}$ .

Let us estimate  $M(\bar{\alpha}_l|\alpha_{<l})$ . Since  $\alpha$  is computable,  $M(\alpha_{<l}) \stackrel{\times}{\geq} 1$ . We claim that  $M(\alpha_{<l} \bar{\alpha}_l) \stackrel{\times}{\leq} 2^{-K(l)}$ . Actually, consider the set  $\{\alpha_{<l} \bar{\alpha}_l | l > 0\}$ . This set is prefix free and decidable. Therefore  $P(l) = M(\alpha_{<l} \bar{\alpha}_l)$  is an enumerable function with  $\sum_l P(l) \leq 1$ , and the claim follows from the coding theorem. Thus, we have  $M(\bar{\alpha}_l|\alpha_{<l}) \stackrel{\times}{\geq} 2^{-K(l)}$  for any  $l$ . Since  $\mu(\bar{\alpha}_l|\alpha_{<l}) > c$ , we get

$$\begin{aligned} \sum_{b \in \{0,1\}} \mu(b|\alpha_{<l}) \ln \frac{\mu(b|\alpha_{<l})}{M(b|\alpha_{<l})} &\stackrel{\pm}{\geq} \mu(\bar{\alpha}_l|\alpha_{<l}) \ln \frac{c}{2^{-K(l)}} + \min_{p \in [0,1-c]} p \ln \frac{p}{M(\alpha_l|\alpha_{<l})} \\ &\stackrel{\pm}{\geq} cK(l) \ln 2 \end{aligned}$$

□



A constant fudge is generally preferable to a logarithmic one for quantitative and aesthetical reasons. It also often leads to particular insight and/or interesting new complexity variants (which will be the case here). Though most complexity variants coincide within logarithmic accuracy (see [Sch00, Sch02a] for exceptions), they can have very different other properties. For instance, Solomonoff complexity  $KM(x) = -\log_2 M(x)$  is an excellent predictor, but monotone complexity  $Km$  can be exponentially worse and prefix complexity  $K$  fails completely [Hut03c].

**Exponential bounds.** Bayes is often approximated by MAP or MDL. In our context this means approximating  $KM$  by  $Km$  with exponentially worse bounds (in deterministic environments) [Hut03c]. (Intuitively, since an error with Bayes eliminates half of the environments, while MAP/MDL may eliminate only one.) Also for more complex “reinforcement” learning problems, bounds can be  $2^{K(\mu)}$  rather than  $K(\mu)$  due to sparser feedback. For instance, for a sequence  $x_1 x_1 x_1 \dots$  if we do not observe  $x_1$  but only receive a reward if our prediction was correct, then the only way a universal predictor can find  $x_1$  is by trying out all  $|\mathcal{X}|$  possibilities and making (in the worst case)  $|\mathcal{X}| - 1 \stackrel{\times}{\cong} 2^{K(\mu)}$  errors. Posterization allows to boost such gross bounds to useful bounds  $2^{K(\mu|x_1)} = O(1)$ . But in general, additive logarithmic corrections as in Theorem 1 also exponentiate and lead to bounds polynomial in  $l$  which may be quite sizeable. Here the advantage of a constant correction becomes even more apparent [Hut04, Problems 2.6, 3.13, 6.3 and Section 5.3.3].

## 5 More Bounds and New Complexity Measure

Lemma 3 shows that the bound in Theorem 1 is attained for some binary strings. But for other binary strings the bound may be very rough. (Similarly,  $K(x)$  is greater than  $\ell(x)$  infinitely often, but  $K(x) \ll \ell(x)$  for many ‘interesting’  $x$ .) Let us try to find a new bound, which does not depend on  $\ell(x)$ .

First observe that, in contrast to the unconditional case (5),  $K(\mu)$  is not an upper bound (again by Lemma 3). Informally speaking, the reason is that  $M$  can predict the future very badly if the past is not “typical” for the environment (such past  $x$  have low  $\mu$ -probability, therefore in the unconditional case their contribution to the expected loss is small). So, it is natural to bound the loss in terms of randomness deficiency  $d_\mu(x)$ , which is a quantitative measure of “typicalness”.

**Theorem 4.** *For any computable measure  $\mu$  and any  $x, y \in \{0, 1\}^*$  it holds*

$$\log_2 \frac{\mu(y|x)}{M(y|x)} \equiv d_\mu(x) - d_\mu(xy) \stackrel{\pm}{\leq} K(\mu) + K(\lceil d_\mu(x) \rceil).$$

Theorem 4 is a variant of the “deficiency conservation theorem” from [VSU05]. We do not know who was the first to discover this statement and whether it was published (the special case where  $\mu$  is the uniform measure was proved by An. Muchnik as an auxiliary lemma for one of his unpublished results; then A. Shen placed a generalized statement to the (unfinished) book [VSU05]).

Now, our goal is to replace  $K(\mu)$  in the last bound by a conditional complexity of  $\mu$ . Unfortunately, the conventional conditional prefix complexity is not suitable:

**Lemma 5.** *Let  $\mathcal{X} = \{0,1\}$ . There is a constant  $C_0$  such that for any  $l \in \mathbb{N}$ , there are a computable measure  $\mu$  and  $x \in \{0,1\}^l$  such that*

$$K(\mu|x) \leq C_0, \quad d_\mu(x) \leq C_0, \quad \text{and}$$

$$D_{l+1:l+1}(x) \equiv \sum_{b \in \{0,1\}} \mu(b|x) \ln \frac{\mu(b|x)}{M(b|x)} \stackrel{\pm}{\geq} K(l) \ln 2.$$

**Proof.** For  $l \in \mathbb{N}$ , define a deterministic measure  $\mu_l$  such that  $\mu_l$  is equal to 1 on the prefixes of  $0^l 1^\infty$  and is equal to 0 otherwise.

Let  $x = 0^l$ . Then  $\mu_l(x) = 1$ ,  $\mu_l(x0) = 0$ ,  $\mu_l(x1) = 1$ . Also  $1 \geq M(x) \geq M(x0) \geq M(0^\infty) \stackrel{\pm}{\geq} 1$  and (as in the proof of Lemma 3)  $M(x1) \stackrel{\pm}{\leq} 2^{-K(l)}$ . Trivially,  $d_{\mu_l}(x) = \log_2 M(x) \stackrel{\pm}{\geq} 1$ , and  $K(\mu_l|x) \stackrel{\pm}{\leq} K(\mu_l|l) \stackrel{\pm}{\leq} 0$ . Thus,  $K(\mu_l|x)$  and  $d_{\mu_l}(x)$  are bounded by a constant  $C_0$  independent of  $l$ . On the other hand,  $\sum_{b \in \{0,1\}} \mu(b|x) \ln \frac{\mu(b|x)}{M(b|x)} = \ln \frac{1}{M(1|x)} \stackrel{\pm}{\geq} K(l) \ln 2$ . (One can obtain the same result also for non-deterministic  $\mu$ , for example, taking  $\mu_l$  mixed with the uniform measure.)  $\square$

Informally speaking, in Lemma 5 we exploit the fact that  $K(y|x)$  can use the information about the length of the condition  $x$ . Hence  $K(y|x)$  can be small for a certain  $x$  and is large for some (actually almost all) prolongations of  $x$ . But in our case of sequence prediction, the length of  $x$  grows taking all intermediate values and cannot contain any relevant information. Thus we need a new kind of conditional complexity.

Consider a Turing machine  $T$  with two input tapes. Inputs are provided without delimiters, so the size of input is defined by the machine itself. Let us call such a machine *twice prefix*. We write that  $T(x,y) = z$  if machine  $T$ , given a sequence beginning with  $x$  on the first tape and a sequence beginning with  $y$  on the second tape, halts after reading exactly  $x$  and  $y$  and prints  $z$  to the output tape. (Obviously, if  $T(x,y) = z$ , then the computation does not depend on the contents of the input tapes after  $x$  and  $y$ .) We define  $C_T(y|x) := \min\{\ell(p) \mid \exists k \leq \ell(x) : T(p, x_{1:k}) = y\}$ . Clearly,  $C_T(y|x)$  is an enumerable from above function of  $T$ ,  $x$ , and  $y$ . Using a standard argument [LV97], one can show that there exists an optimal twice prefix machine  $U$  in the sense that for any twice prefix machine  $T$  we have  $C_U(y|x) \stackrel{\pm}{\leq} C_T(y|x)$ .

**Definition 6.** *Complexity monotone in conditions* is defined for some fixed optimal twice prefix machine  $U$  as

$$K_*(y|x^*) := C_U(y|x) = \min\{\ell(p) \mid \exists k \leq \ell(x) : U(p, x_{1:k}) = y\}.$$

Here  $*$  in  $x^*$  is a syntactical part of the complexity notation, though one may think of  $K_*(y|x^*)$  as of the minimal length of a program that produces  $y$  given any  $z = x^*$ .

**Theorem 7.** *For any computable measure  $\mu$  and any  $x, y \in \mathcal{X}^*$  it holds*

$$\log_2 \frac{\mu(y|x)}{M(y|x)} \stackrel{\pm}{\leq} K_*(\mu|x^*) + K(\lceil d_\mu(x) \rceil).$$

*Note.* One can get a slightly stronger variants of Theorems 1 and 7 by replacing the complexity of a standard code of  $\mu$  by more sophisticated values. First, in any effective encoding there are many codes for every  $\mu$ , and in all the upper bounds (including Solomonoff's one) one can take the minimum of the complexities of all the codes for  $\mu$ . Moreover, in Theorem 1 it is sufficient to take the complexity of  $\mu_x = \mu(\cdot|x)$  (and it is sufficient that  $\mu_x$  is enumerable, while  $\mu$  can be incomputable). For Theorem 7 one can prove a similar strengthening: The complexity of  $\mu$  is replaced by the complexity of any computable function that is equal to  $\mu$  on all prefixes and prolongations of  $x$ .

To demonstrate the usefulness of the new bound, let us again consider some deterministic environment  $\mu \hat{=} \alpha$ . For  $\mathcal{X} = \{0,1\}$  and  $\alpha = x^\infty$  with  $x = 0^n 1$ , Theorem 1 gives the bound  $K(\mu|n) + K(n) \stackrel{\pm}{\leq} K(n)$ . Consider the new bound  $K_*(\mu|x^*) + K(\lceil d_\mu(x) \rceil)$ . Since  $\mu$  is deterministic, we have  $d_\mu(x) = \log_2 M(x) \stackrel{\pm}{\leq} -K(n)$ , and  $K(\lceil d_\mu(x) \rceil) \stackrel{\pm}{\leq} K(K(n))$ . To estimate  $K_*(\mu|x^*)$ , let us consider a machine  $T$  that reads only its second tape and outputs the number of 0s before the first 1. Clearly,  $C_T(n|x) = 0$ , hence  $K_*(\mu|x^*) \stackrel{\pm}{\leq} 0$ . Finally,  $K_*(\mu|x^*) + K(\lceil d_\mu(x) \rceil) \stackrel{\pm}{\leq} K(K(n))$ , which is much smaller than  $K(n)$ .

## 6 Properties of the New Complexity

The above definition of  $K_*$  is based on computations of some Turing machine. Such definitions are quite visual, but are often not convenient for formal proofs. We will give an alternative definition in terms of enumerable sets (see [US96] for definitions of unconditional complexities in this style), which summarizes the properties we actually need for the proof of Theorem 7.

An enumerable set  $E$  of triples of strings is called  $K_*$ -correct if it satisfies the following requirements:

1. if  $\langle p, x, y_1 \rangle \in E$  and  $\langle p, x, y_2 \rangle \in E$ , then  $y_1 = y_2$ ;
2. if  $\langle p, x, y \rangle \in E$ , then  $\langle p', x', y \rangle \in E$  for all  $p'$  being prolongations of  $p$  and all  $x'$  being prolongations of  $x$ ;
3. if  $\langle p, x', y \rangle \in E$  and  $\langle p', x, y \rangle \in E$ , and  $p$  is a prefix of  $p'$  and  $x$  is a prefix of  $x'$ , then  $\langle p, x, y \rangle \in E$ .

A complexity of  $y$  under a condition  $x$  w.r.t. a set  $E$  is  $C_E(y|x) = \min\{\ell(p) \mid \langle p, x, y \rangle \in E\}$ . A  $K_*$ -correct set  $E$  is called *optimal* if  $C_E(y|x) \stackrel{\pm}{\leq} C_{E'}(y|x)$  for any  $K_*$ -correct set  $E'$ . One can easily construct an enumeration of all  $K_*$ -correct sets, and an optimal set exists by the standard argument.

It is easy to see that a twice prefix Turing machine  $T$  can be transformed to a set  $E$  such that  $C_T(y|x) = C_E(y|x)$ . The set  $E$  is constructed as follows:  $T$  is run on all possible inputs, and if  $T(p, x) = y$ , then pairs  $\langle p', x', y \rangle$  are added to  $E$  for all  $p'$  being prolongations of  $p$  and all  $x'$  being prolongations of  $x$ . Evidently,  $E$  is enumerable, and the second requirement of  $K_*$ -correctness is satisfied. To verify the

other requirements, let us consider arbitrary  $\langle p'_1, x'_1, y_1 \rangle \in E$  and  $\langle p'_2, x'_2, y_2 \rangle \in E$  such that  $p'_1$  and  $p'_2$ ,  $x'_1$  and  $x'_2$  are comparable (one is a prefix of the other). Then, by construction of  $E$ , we have  $T(p_1, x_1) = y_1$  and  $T(p_2, x_2) = y_2$ , and  $p_1$  and  $p_2$ ,  $x_1$  and  $x_2$  are comparable too. Since replacing the unused part of the inputs does not affect the running of the machine  $T$  and comparable words have a common prolongation, we get  $p_1 = p_2$ ,  $x_1 = x_2$ , and  $y_1 = y_2$ . Thus  $E$  is a  $K_*$ -correct set.

The transformation in the other direction is impossible in some cases: the set  $E = \{\langle 0^{h(n)}p, 0^n 1q, 0 \rangle \mid n \in \mathbb{N}, p, q \in \{0, 1\}^*\}$ , where  $h(n)$  is 0 if the  $n$ -th Turing machine halts and 1 otherwise, is  $K_*$ -correct, but does not have a corresponding machine  $T$ : using such a machine one could solve the halting problem. However, we conjecture that for every set  $E$  there exists a machine  $T$  such that  $C_T(x|y) \stackrel{\pm}{\leq} C_E(x|y)$ .

Probably, the requirements on  $E$  can be even weaker, namely, the third requirement can be superfluous. Let us notice that the first requirement of  $K_*$ -correctness allows us to consider the set  $E$  as a partial computable function:  $E(p, x) = y$  iff  $\langle p, x, y \rangle \in E$ . The second requirement says that  $E$  becomes a continuous function if we take the topology of prolongations (any neighborhood of  $\langle p, x \rangle$  contains the cone  $\{\langle p^*, x^* \rangle\}$ ) on the arguments and the discrete topology ( $\{y\}$  is a neighborhood of  $y$ ) on values. It is known (see [US96] for references) that different complexities (plain, prefix, decision) can be naturally defined in a similar “topological” fashion. We conjecture the same is true in our case: an optimal enumerable set satisfying the requirements (1) and (2) (obviously, it exists) specifies the same complexity (up to an additive constant) as an optimal twice prefix machine.

It follows immediately from the definition(s) that  $K_*(y|x^*)$  is monotone as a function of  $x$ :  $K_*(y|xz^*) \leq K_*(y|x^*)$  for all  $x, y, z$ .

The following lemma provides bounds for  $K_*(x|y^*)$  in terms of prefix complexity  $K$ . The lemma holds for all our definitions of  $K_*(x|y^*)$ .

**Lemma 8.** *For any  $x, y \in \mathcal{X}^*$  it holds*

$$K(x|y) \stackrel{\pm}{\leq} K_*(x|y^*) \stackrel{\pm}{\leq} \min_{l \leq \ell(y)} \{K(x|y_{1:l}) + K(l)\} \stackrel{\pm}{\leq} K(x).$$

*In general, none of the bounds is equal to  $K_*(x|y^*)$  even within  $o(K(x))$  term, but they are attained for certain  $y$ : For every  $x$  there is a  $y$  such that*

$$K(x|y) \stackrel{\pm}{=} 0 \quad \text{and} \quad K_*(x|y^*) \stackrel{\pm}{=} K(x) \stackrel{\pm}{=} \min_{l \leq \ell(y)} \{K(x|y_{1:l}) + K(l)\},$$

*and for every  $x$  there is a  $y$  such that*

$$K(x|y) \stackrel{\pm}{=} K_*(x|y^*) \stackrel{\pm}{=} 0 \quad \text{and} \quad K(x) \stackrel{\pm}{\leq} \min_{l \leq \ell(y)} \{K(x|y_{1:l}) + K(l)\}.$$

**Corollary 9.** *The future deviation of  $M_t$  from  $\mu_t$  is bounded by*

$$\sum_{t=l+1}^{\infty} \mathbf{E}[s_t | \omega_{1:l}] \stackrel{\pm}{\leq} [\min_{i \leq l} \{K(\mu | \omega_{1:i}) + K(i)\} + K(d_\mu(\omega_{1:l}))] \ln 2.$$

Let us note that if  $\omega$  is  $\mu$ -random, then  $K(d_\mu(\omega_{1:l})) \stackrel{\pm}{\leq} K(d_\mu(\omega_{1:\infty})) + K(K(\mu))$ , and therefore we get the bound, which does not increase with  $l$ , in contrast to the bound (i) in Corollary 2.

## 7 Proof of Theorem 7

The plan is to get a statement of the form  $2^d \mu(y) \stackrel{\times}{\leq} M(y)$ , where  $d \approx d_\mu(x) = \log_2 \frac{M(x)}{\mu(x)}$ . To this end, we define a new semimeasure  $\nu$ : we take the set  $S = \{z \mid d_\mu(z) > d\}$  and put  $\nu$  to be  $2^d \mu$  on prolongations of  $z \in S$ ; this is possible since  $S$  has  $\mu$ -measure  $2^{-d}$ . Then we have  $\nu(z) \leq C \cdot M(z)$  by universality of  $M$ . However, the constant  $C$  depends on  $\mu$  and also on  $d$ . To make the dependence explicit, we repeat the above construction for all numbers  $d$  and all semimeasures  $\mu^T$ , obtaining semimeasures  $\nu_{d,T}$ , and take  $\nu = \sum 2^{-K(d)} \cdot 2^{-K(T)} \nu_{d,T}$ . This construction would give us the term  $K(\mu)$  in the right-hand side of Theorem 7. To get  $K_*(\mu|x^*)$ , we need a more complicated strategy: instead of a sum of semimeasures  $\nu_{d,T}$ , for every fixed  $d$  we sum “pieces” of  $\nu_{d,T}$  at each point  $z$ , with coefficients depending on  $z$  and  $T$ .

Now proceed with the formal proof. Let  $\{\mu^T\}_{T \in \mathcal{N}}$  be any (effective) enumeration of all enumerable semimeasures. For any integer  $d$  and any  $T$ , put

$$S_{d,T} := \{z \mid \sum_{v \in \mathcal{X}^{\ell(z)} \setminus \{z\}} \mu^T(v) + 2^{-d} M(z) > 1\}.$$

The set  $S_{d,T}$  is enumerable given  $d$  and  $T$ .

Let  $E$  be the optimal  $K_*$ -correct set (satisfying all three requirements),  $E(p,z)$  is the corresponding partial computable function. For any  $z \in \mathcal{X}^*$  and  $T$ , put

$$\lambda(z, T) := \max\{2^{-\ell(p)} \mid \exists k \leq \ell(z) : z_{1:k} \in S_{d,T} \text{ and } E(p, z_{1:k}) = T\}$$

(if there is no such  $p$ , then  $\lambda(z, T) = 0$ ). Put

$$\tilde{\nu}_d(z) := \sum_T \lambda(z, T) \cdot 2^d \mu^T(z).$$

Obviously, this value is enumerable. It is not a semimeasure, but it has the following property (we omit the proof).

**Claim 10.** *For any prefix-free set  $A$ ,*

$$\sum_{z \in A} \tilde{\nu}_d(z) \leq 1.$$

This implies that there exists an enumerable semimeasure  $\nu_d$  such that  $\nu_d(z) \geq \tilde{\nu}_d(z)$  for all  $z$ . Actually, to enumerate  $\nu_d$ , one enumerates  $\tilde{\nu}_d(z)$  for all  $z$  and at each step sets the current value of  $\nu_d(z)$  to the maximum of the current values of  $\tilde{\nu}_d(z)$  and  $\sum_{u \in \mathcal{X}} \nu_d(zu)$ . Trivially, this provides  $\nu_d(z) \geq \sum_{u \in \mathcal{X}} \nu_d(zu)$ . To show that  $\nu_d(\epsilon) \leq 1$ , let us note that at any step of enumeration the current value of  $\nu_d(\epsilon)$  is the sum of current values  $\tilde{\nu}_d(z)$  over some prefix-free set, and thus is bounded by 1. Put

$$\nu(z) := \sum_d 2^{-K(d)} \nu_d(z).$$

Clearly,  $\nu$  is an enumerable semimeasure, thus  $\nu(z) \stackrel{\times}{\leq} M(z)$ . Let  $\mu$  be an arbitrary computable measure, and  $x, y \in \mathcal{X}^*$ . Let  $p \in \{0,1\}^*$  be a string such that  $K_*(\mu|x^*) =$

$\ell(p)$ ,  $E(p,x)=T$ , and  $\mu=\mu^T$ . Put  $d=\lceil d_\mu(x)\rceil-1$ , i.e.,  $d_\mu(x)-1\leq d<d_\mu(x)$ . Hence  $\mu(x)<2^{-d}M(x)$ . Since  $\mu=\mu^T$  is a measure, we have  $\sum_{v\in\mathcal{X}^{\ell(x)}}\mu^T(v)=1$ , and therefore  $x\in S_{d,T}$ . By definition,  $\lambda(xy,T)\geq 2^{-\ell(p)}$ , thus  $\tilde{\nu}_d(xy)\geq 2^{-\ell(p)}2^d\mu(xy)$ , and

$$2^{-K(d)}2^{-\ell(p)}2^d\mu(xy)\leq\nu(xy)\stackrel{\times}{\leq}M(xy).$$

After trivial transformations we get

$$\log_2\frac{\mu(y|x)}{M(y|x)}\stackrel{\pm}{\leq}K_*(\mu|x^*)+K(d),$$

which completes the proof of Theorem 7.

## 8 Discussion

**Conclusion.** We evaluated the quality of predicting a stochastic sequence at an intermediate time, when some beginning of the sequence has been already observed, estimating the future loss of the universal Solomonoff predictor  $M$ . We proved general upper bounds for the discrepancy between conditional values of the predictor  $M$  and the true environment  $\mu$ , and demonstrated a kind of tightness for these bounds. One of the bounds is based on a new variant of conditional algorithmic complexity  $K_*$ , which has interesting properties in its own. In contrast to standard prefix complexity  $K$ ,  $K_*$  is a monotone function of conditions:  $K_*(y|xyz^*)\leq K_*(y|x^*)$ .

**General Bayesian posterior bounds.** A natural question is whether posterior bounds for general Bayes mixtures based on general  $\mathcal{M}\ni\mu$  could also be derived. From the (obvious) posterior representation  $\xi(y|x)=\sum_{\nu\in\mathcal{M}}w_\nu(x)\nu(y|x)\geq w_\mu(x)\mu(y|x)$ , where  $w_\nu(x):=w_\nu\frac{\nu(x)}{\xi(x)}$  is the posterior belief in  $\nu$  after observing  $x$ , the bound  $D_{l:\infty}\leq\ln w_\mu(\omega_{<l})^{-1}$  immediately follows. Strangely enough, for  $\mathcal{M}=\mathcal{M}_U$ ,  $\log_2 w_\nu^{-1}:=K(\nu)$  does *not* imply  $\log_2 w_\mu(x)^{-1}=K(\mu|x)$ , not even within logarithmic accuracy, so it was essential to consider  $D_{l:\infty}$ . It would be interesting to derive bounds on  $D_{l:\infty}$  or  $\ln w_\mu(x)^{-1}$  for general  $\mathcal{M}$  similar to the ones derived here for  $\mathcal{M}=\mathcal{M}_U$ .

**Online classification.** All considered distributions  $\rho(x)$  (in particular  $\xi$ ,  $M$ , and  $\mu$ ), may be replaced everywhere by distributions  $\rho(x|z)$  additionally conditioned on some  $z$ . The  $z$ -conditions cause nowhere problems as they can essentially be thought of as fixed (or as oracles or spectators). An (i.i.d.) classification problem is a typical example: At time  $t$  one arranges an experiment  $z_t$  (or observes data  $z_t$ ), then tries to make a prediction, and finally observes the true outcome  $x_t$  with probability  $\mu(x_t|z_t)$ . In this case  $\mathcal{M}=\{\nu(x_{1:n}|z_{1:n})=\nu(x_1|z_1)\cdots\nu(x_n|z_n)\}$ . (Note that  $\xi$  is not i.i.d.) Solomonoff's bound  $K(\mu)\ln 2$  (5) holds unchanged. Compared to the sequence prediction case we have extra information  $z$ , so we may wonder whether some improved bound  $K(\mu|z)$  or so, holds. For a fixed  $z$  this can be achieved by also replacing  $2^{-K(\mu)}$  in (2) by  $2^{-K(\mu|z)}$ . But if at time  $t$  only  $z_{1:t}$  is known like in the classification example, this leads to difficulties ( $\xi$  is no longer a (semi)measure, which

sometimes can be corrected [PH04]). Alternatively we could keep definition (2) but apply it to the (chronologically correctly ordered) sequence  $z_1x_1z_2x_2\dots$ , condition to  $z_{1:t}$ , and try to derive improved bounds.

**More open problems.** Since  $D_{1:\infty}$  is finite, one may expect that the tails  $D_{l:\infty}$  tend to 0 as  $l \rightarrow \infty$ . However, as Lemma 3 implies, this holds only with probability 1: for some special  $\alpha$  we have even  $D_{l:\infty}(\alpha_{<l}) \stackrel{\pm}{\geq} \frac{1}{3}K(l) \xrightarrow{l \rightarrow \infty} \infty$ . It would be very interesting to find a wide class of  $\alpha$  such that  $D_{l:\infty}(\alpha_{<l}) \rightarrow 0$ . The natural conjecture is that one should take  $\mu$ -random  $\alpha$ . Another (probably, closely related) task is to study the asymptotic behavior of  $K_*(\mu|\alpha_{<l}^*)$ . It is natural to expect that  $K_*(\mu|\alpha_{<l}^*)$  is bounded by an absolute constant (independent of  $\mu$ ) for “most”  $\alpha$  and for sufficiently large  $l$ . Finally, (dis)proving equality of the various definitions of  $K_*$  we gave, would be useful.

## References

- [CV05] R. Cilibrasi and P. M. B. Vitányi. Clustering by compression. *IEEE Trans. Information Theory*, 51(4):1523–1545, 2005.
- [HM04] M. Hutter and An. A. Muchnik. Universal convergence of semimeasures on individual random sequences. In *Proc. 15th International Conf. on Algorithmic Learning Theory (ALT’04)*, volume 3244 of *LNAI*, pages 234–248, Padova, 2004. Springer, Berlin.
- [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’03)*, volume 2777 of *LNAI*, pages 506–521, Berlin, 2003. Springer.
- [Hut04] M. Hutter. *Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability*. Springer, Berlin, 2004. 300 pages, <http://www.idsia.ch/~marcus/ai/uaibook.htm>.
- [LV97] M. Li and P. M. B. Vitányi. *An introduction to Kolmogorov complexity and its applications*. Springer, 2nd edition, 1997.
- [PH04] J. Poland and M. Hutter. Convergence of discrete MDL for sequential prediction. In *Proc. 17th Annual Conf. on Learning Theory (COLT’04)*, volume 3120 of *LNAI*, pages 300–314, Banff, 2004. Springer, Berlin.

- [Sch00] J. Schmidhuber. Algorithmic theories of everything. Report IDSIA-20-00, quant-ph/0011122, IDSIA, Manno (Lugano), Switzerland, 2000.
- [Sch02a] J. Schmidhuber. Hierarchies of generalized Kolmogorov complexities and nonenumerable universal measures computable in the limit. *International Journal of Foundations of Computer Science*, 13(4):587–612, 2002.
- [Sch02b] J. Schmidhuber. The Speed Prior: a new simplicity measure yielding near-optimal computable predictions. In *Proc. 15th Annual Conference on Computational Learning Theory (COLT 2002)*, Lecture Notes in Artificial Intelligence, pages 216–228, Sydney, Australia, July 2002. Springer.
- [Sol64] R. J. Solomonoff. A formal theory of inductive inference: Part 1 and 2. *Inform. Control*, 7:1–22, 224–254, 1964.
- [Sol78] R. J. Solomonoff. Complexity-based induction systems: comparisons and convergence theorems. *IEEE Trans. Information Theory*, IT-24:422–432, 1978.
- [US96] V. A. Uspensky and A. Shen. Relations Between Varieties of Kolmogorov Complexities. *Math. Systems Theory*, 29:271–292, 1996.
- [VSU05] N. K. Vereshchagin, A. Shen, and V. A. Uspensky. Lecture Notes on Kolmogorov Complexity. Unpublished, <http://lpcs.math.msu.su/~ver/kolm-book>, 2005.
- [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.