



**HAL**  
open science

## Computable priors sharpened into Occam's razors

David R. Bickel

► **To cite this version:**

| David R. Bickel. Computable priors sharpened into Occam's razors. 2016. hal-01423673v2

**HAL Id: hal-01423673**

**<https://hal.science/hal-01423673v2>**

Preprint submitted on 23 Jan 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Copyright

# Computable priors sharpened into Occam's razors

David R. Bickel

January 23, 2017

Ottawa Institute of Systems Biology

Department of Biochemistry, Microbiology, and Immunology

Department of Mathematics and Statistics

University of Ottawa

451 Smyth Road

Ottawa, Ontario, K1H 8M5

+01 (613) 562-5800, ext. 8670

dbickel@uottawa.ca

## Abstract

The posterior probabilities available under standard Bayesian statistics are computable, apply to small samples, and coherently incorporate previous information. Modifying their priors according to results from algorithmic information theory adds the advantage of implementing Occam's razor, giving simpler distributions of data higher prior probabilities.

**Keywords:** algorithmic probability; Bayesian inference; entropy rate; Kolmogorov complexity; universal prior

Two general formalisms for data analysis based on Bayes's theorem are standard Bayesian statistics and Solomonoff's approach to universal distributions under algorithmic information theory. The posterior probabilities available under the former are computable and explicitly represent previous knowledge in terms of probability distributions. While lacking those advantages, Solomonoff's use of universal distributions implements Occam's razor by assigning more prior probability to sequences of symbols that are simpler in the sense of having lower Kolmogorov complexity (see Solomonoff, 1978, 2008).

In this note, the two approaches are combined to make their advantages simultaneously applicable to data analysis, inference, and decision. The complexity of each stochastic process under consideration will determine its prior probability for Bayesian inference and decision making consistent with Occam's razor. That prior probability of a process will be defined as the limit of its probability that an observation coincides with a symbol distributed according to a universal prior under algorithmic information theory.

To make that precise, consider a stationary process  $X_{\theta,1}, X_{\theta,2}, \dots$  such that the law of  $X_{\theta,t}$  is a probability mass function  $f_{\theta,t}$  on a finite alphabet  $\mathcal{X}$  for all  $t = 1, 2, \dots$  and for each parameter value  $\theta$  in a set  $\Theta$ . *Strings*, finite sequences of  $\tau$  symbols (Nies, 2012, §1.2), are written as  $X_{\theta}^{\tau} = (X_{\theta,1}, \dots, X_{\theta,\tau})$ . The length of any string  $x$  is  $\ell(x)$ ; thus,  $\ell(X_{\theta}^{\tau}) = \tau$ . Throughout,  $\log = \log_2$ . Convergence in probability and weak convergence as  $\tau \rightarrow \infty$  are both denoted by  $\xrightarrow{\text{weak}}$  (see Billingsley, 1999, p. 27).

Let  $\vartheta$  denote a random variable with distribution  $\pi_0$ , a probability measure on the measurable space  $(\Theta, \mathfrak{H})$ , where  $\mathfrak{H}$  is a  $\sigma$ -field of subsets of  $\Theta$ . For any  $\mathcal{H} \in \mathfrak{H}$ ,  $\pi_0(\mathcal{H})$  is then  $P(\vartheta \in \mathcal{H})$ , the probability of the hypothesis that the process  $X_{\theta,1}, X_{\theta,2}, \dots$  is in  $\{(X_{\theta,1}, X_{\theta,2}, \dots) : \theta \in \mathcal{H}\}$ . Since  $\pi_0$  enables but does not depend on the complexity considerations of Occam's razor, it is called the *blunted prior distribution*. It is used below to incorporate any application-specific information into the *universal prior*  $Q$ , defined such that  $Q(x)$  is the algorithmic probability that a given prefix-free machine generates output  $x$  from a program consisting of independent and equally probable symbols (Downey, 2008; Nies, 2012, §2.2; cf. Hutter, 2006, §2.4.1).

As a discrete semimeasure, the total probability mass of  $Q$  may be less than 1 (Downey and Hirschfeldt, 2010, §3.9). For that reason, normalization is required to transform  $Q$  to a probability mass function before applying standard probability theory. According to the *universal probability mass function*  $Q_{\text{norm}}$ ,

the conditional probability of a symbol following  $\tau$  previous symbols based on  $Q$  is

$$Q_{\text{norm}}(x_{\tau+1}|x^\tau) = \frac{Q_{\text{norm}}(x^{\tau+1})}{Q_{\text{norm}}(x^\tau)} \propto \frac{Q(x^{\tau+1})}{Q(x^\tau)}, \quad (1)$$

normalized such that  $\sum_{y \in \mathcal{X}} Q_{\text{norm}}(y|x^\tau) = 1$ , as required for  $Q_{\text{norm}}(\bullet|x^\tau)$  to be a probability mass function, for all  $x^\tau \in \mathcal{X}^\tau$  (cf. Li and Vitányi, 2008, §4.5.3). While  $Q_{\text{norm}}$  is incomputable and was originally intended for making predictions without any blunted  $\pi_0$  or other explicit reliance on prior knowledge (see Rathmanner and Hutter, 2011, §10.1), the leverage of such a  $\pi_0$  leads to a computable posterior distribution suitable for predictions and other Bayes actions, decisions that minimize posterior expected loss.

Let  $U_{\theta, \tau+1}$  be a random variable distributed according to the random probability mass function  $Q_{\text{norm}}(\bullet|X_\theta^\tau)$  for each  $\theta \in \Theta$ . (The randomness of  $Q_{\text{norm}}(\bullet|X_\theta^\tau)$  comes from that of  $X_\theta^\tau$  and is meant in the sense of probability theory, not in the sense in which an individual string can be random.) For any  $\mathcal{H} \in \mathfrak{H}$  and  $t = 1, 2, \dots$ , the  $\tau$ -whetted probability  $\pi_\tau(\mathcal{H})$  and the whetted probability  $\pi(\mathcal{H})$  that  $\vartheta \in \mathcal{H}$  are defined by

$$\pi_\tau(\mathcal{H}) = P(\vartheta \in \mathcal{H} | X_{\vartheta, \tau+1} = U_{\vartheta, \tau+1}) \quad (2)$$

$$\pi(\mathcal{H}) = \frac{\pi_0(\mathcal{H}) \int_{\mathcal{H}} d\pi_0(\theta) 2^{-H(\theta)}}{\int d\pi_0(\theta) 2^{-H(\theta)}}, \quad (3)$$

where the function  $H$  on  $\Theta$  gives the entropy rate  $H(\theta) = \lim_{\tau \rightarrow \infty} H(X_{\theta,1}, \dots, X_{\theta,\tau}) / \tau$  on the basis of the Shannon entropy  $H(X_{\theta,1}, \dots, X_{\theta,\tau})$  for each  $\theta \in \Theta$ . Equations (2)-(3) define two prior probability measures on  $(\Theta, \mathfrak{H})$ , the  $\tau$ -whetted distribution  $\pi_\tau$  and the whetted distribution  $\pi$ . Describing both with the same adjective will be justified by  $\pi_\tau \xrightarrow{\text{weak}} \pi$ .

**Example 1.** Suppose  $\Theta$  is the set of nonnegative integers. For any  $\theta \in \Theta$  and  $t = 1, 2, \dots$ , the  $\tau$ -whetted probability  $\pi_\tau(\theta)$  and the whetted probability  $\pi(\theta)$  that the process is  $X_{\theta,1}, X_{\theta,2}, \dots$  are, with  $\pi_0(\theta) = \pi_0(\{\theta\})$ ,

$$\begin{aligned} \pi_\tau(\theta) &= \pi_\tau(\{\theta\}) = P(\vartheta = \theta | X_{\vartheta, \tau+1} = U_{\vartheta, \tau+1}) \\ \pi(\theta) &= \pi(\{\theta\}) = \frac{\pi_0(\theta) 2^{-H(\theta)}}{\sum_{i \in \Theta} \pi_0(i) 2^{-H(i)}}. \end{aligned} \quad (4)$$

**Example 2.** In the finite- $\Theta$  setting, suppose  $X_{\theta,1}, X_{\theta,2}, \dots$  are mutually independent,  $\Theta = \mathcal{X} = \{1, 2, \dots, m\}$ ,  $\pi_0(\theta) = 1/m$ , and  $f_{\theta,t}(x) = 1/\theta$  if  $x \leq \theta$  but  $f_{\theta,t}(x) = 0$  if  $x > \theta$  for all  $\theta, x = 1, 2, \dots, m$  and  $t = 1, 2, \dots$

The entropy rate  $H(\theta)$  is the entropy  $-\sum_{x=1}^m f_{\theta,t}(x) \log f_{\theta,t}(x) = \log \theta$  under independence. Thus, according to equation (4), the whetted distribution is given by  $\pi(\theta) \propto m^{-1} 2^{-\log \theta} = m^{-1} \theta^{-1}$ .

Under ergodicity, the  $\tau$ -whetted distribution converges to the whetted distribution, which is computable since  $\pi_0$  and  $f_{\theta,t}$  are computable.

**Theorem 1.** *If the process  $X_{\theta,1}, X_{\theta,2}, \dots$  is ergodic for every  $\theta \in \Theta$ , then  $\pi_\tau \xrightarrow{\text{weak}} \pi$ .*

*Proof.* Let  $\mathcal{H}$  denote any continuity set in  $\mathfrak{S}$ . For any  $t = 1, 2, \dots$ ,

$$\pi_\tau(\mathcal{H}) = \frac{P(\vartheta \in \mathcal{H}) P(U_{\vartheta, \tau+1} = X_{\vartheta, \tau+1} | \vartheta \in \mathcal{H})}{P(U_{\vartheta, \tau+1} = X_{\vartheta, \tau+1})} = \frac{\pi_0(\mathcal{H}) \int_{\mathcal{H}} d\pi_0(\theta | \mathcal{H}) P(U_{\theta, \tau+1} = X_{\theta, \tau+1})}{P(U_{\vartheta, \tau+1} = X_{\vartheta, \tau+1})}. \quad (5)$$

The key factor is  $\int_{\mathcal{H}} d\pi_0(\theta | \mathcal{H}) P(U_{\theta, \tau+1} = X_{\theta, \tau+1})$ , in which

$$\begin{aligned} P(U_{\theta, \tau+1} = X_{\theta, \tau+1}) &= \sum_{x \in \mathcal{X}} P(X_{\theta, \tau+1} = x) P(U_{\theta, \tau+1} = X_{\theta, \tau+1} | X_{\theta, \tau+1} = x) \\ &= \sum_{x \in \mathcal{X}} f_{\theta, \tau+1}(x) P(U_{\theta, \tau+1} = x | X_{\theta, \tau+1} = x) \\ &= \sum_{x \in \mathcal{X}} f_{\theta, \tau+1}(x) E(Q_{\text{norm}}(x | X_\theta^\tau) | X_{\theta, \tau+1} = x) \\ &= E(Q_{\text{norm}}(X_{\theta, \tau+1} | X_\theta^\tau)) = E\left(\frac{Q_{\text{norm}}(X_{\theta, \tau+1})}{Q_{\text{norm}}(X_\theta^\tau)}\right) \end{aligned} \quad (6)$$

for all  $\theta \in \mathcal{H}$ , with the last step from equation (1). Since  $X_{\theta,1}, X_{\theta,2}, \dots$  is an ergodic process,  $K(X_\theta^\tau | \tau) / \tau \xrightarrow{\text{weak}} H(\theta)$ , where  $K(x | \tau)$  is the prefix Kolmogorov complexity of a sequence  $x$  conditional on  $\ell(x) = \tau$  (Horibe, 2003). That implies  $K(X_\theta^\tau) / \tau \xrightarrow{\text{weak}} H(\theta)$  (cf. Li and Vitányi, 2008, p. 605), for there are constants  $c_1, c_2 > 0$  that satisfy

$$K(x^\tau | \tau) - c_1 \leq K(x^\tau) \leq K(x^\tau | \tau) + 2 \log \tau + c_2 \quad \forall x^\tau \in \mathcal{X}^\tau$$

with the lower bound from  $K(x^\tau | \tau) \leq K(x^\tau) + c_1$ , following Li and Vitányi (2008, p. 119), and the upper bound from Cover and Thomas (2006, Theorem 14.2.3); cf. Li and Vitányi (2008, p. 204). Writing  $o_P(U_\tau)$  for any random sequence such that  $o_P(U_\tau) / U_\tau \xrightarrow{\text{weak}} 0$  given a random sequence  $U_\tau$  (Serfling, 1980, §1.2.5),

$$\frac{K(X_\theta^\tau)}{\tau} = \frac{K(X_\theta^\tau | \tau)}{\tau} + o_P(1) = H(\theta) + o_P(1). \quad (7)$$

The coding theorem (e.g., Downey, 2008, Theorem 2.8; Nies, 2012, Theorem 2.2.25) entails that there are a  $c_3 > 0$  and a  $c_4 > 0$  satisfying

$$K(x^\tau) - c_3 \leq -\log Q(x^\tau) \leq K(x^\tau) + c_4.$$

That yields  $(-\log Q(X_\theta^\tau) - K(X_\theta^\tau)) / \tau \xrightarrow{\text{weak}} 0$  and thus, with equation (7),

$$\frac{-\log Q(X_\theta^\tau)}{\tau} = \frac{K(X_\theta^\tau)}{\tau} + o_P(1) = H(\theta) + o_P(1).$$

$$\begin{aligned} \therefore -\log Q(X_\theta^{\tau+1}) - (-\log Q(X_\theta^\tau)) &= (\tau+1) \left( \frac{K(X_\theta^{\tau+1})}{\tau+1} + o_P(1) \right) - \tau \left( \frac{K(X_\theta^\tau)}{\tau} + o_P(1) \right) \\ &= (\tau+1)(H(\theta) + o_P(1)) - \tau(H(\theta) + o_P(1)) \xrightarrow{\text{weak}} H(\theta). \end{aligned}$$

The continuous mapping theorem then yields

$$\frac{Q(X_\theta^{\tau+1})}{Q(X_\theta^\tau)} = 2^{-(\log Q(X_\theta^{\tau+1}) - \log Q(X_\theta^\tau))} \xrightarrow{\text{weak}} 2^{-H(\theta)}.$$

$$\therefore \lim_{\tau \rightarrow \infty} P(U_{\theta, \tau+1} = X_{\theta, \tau+1}) \propto 2^{-H(\theta)}$$

according to equation (6) since  $\exp(-(Q(X_\theta^{\tau+1}) - Q(X_\theta^\tau)))$  is bounded. From equation (5),  $\pi_\tau(\mathcal{H})$  converges to a quantity proportional to  $\pi_0(\mathcal{H}) \int_{\mathcal{H}} d\pi_0(\theta|\mathcal{H}) \exp(-H(\theta))$ . Since  $\int d\pi_\tau(\theta) = 1$  for all  $\tau = 1, 2, \dots$ , it follows that  $\pi_\tau(\mathcal{H}) \rightarrow \pi(\mathcal{H})$  for every continuity set  $\mathcal{H}$  in  $\mathfrak{H}$ . The portmanteau theorem makes that equivalent to  $\pi_\tau \xrightarrow{\text{weak}} \pi$ .  $\square$

Since  $\theta \mapsto 2^{-H(\theta)}$  in equation (3) is mathematically equivalent to a likelihood function and yet does not depend on any data, it is an example of a prior likelihood function, the logarithm of which is a “prior support” function as defined by Edwards (1992). As likelihood is defined only up to a multiplicative constant,  $L(\bullet) = c2^{-H(\bullet)}$  for any  $c > 0$  may be called the *whetted likelihood function*. It relates the *whetted* and *blunted probability densities* by  $(d\pi/d\eta)(\theta) \propto L(\theta)(d\pi_0/d\eta)(\theta)$  for any  $\theta \in \Theta$ , where  $\eta$  is any measure that dominates  $\pi_0$  and  $\pi$ . In the sense that the use of the whetted likelihood justifies the frequent claim that the prior probability density should be higher for simpler hypotheses, it resembles the *simplicity postulate*, which

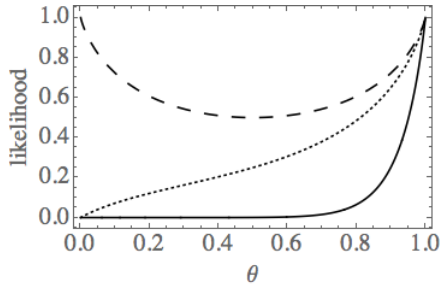


Figure 1: Likelihood functions for the binomial probability as  $\theta$ . The whetted likelihood function is combined with the likelihood function for  $n = 0$  (dashed),  $n = 1$  (dotted), and  $n = 10$  (solid) consecutive successes.

requires laws with fewer free parameters to have higher prior probabilities (Jeffreys, 1948, pp. 100-101).

**Example 3.** Let  $X_{\theta,1}, X_{\theta,2}, \dots$  denote independent Bernoulli variates of unknown probability  $\theta$  of success. By independence, the entropy rate  $H(\theta)$  is the entropy,  $-\theta \log \theta - (1 - \theta) \log (1 - \theta)$ . The resulting whetted likelihood  $L(\theta) = \theta^\theta (1 - \theta)^{(1-\theta)}$  is the dashed curve in Figure 1. The other two plotted curves are the products of  $L(\theta)$  and the likelihood functions of the binomial distribution for  $n = 1$  and  $n = 10$  observations, all successes. (The finite sample size  $n$  should not be confused with  $\tau$ , which goes to infinity.) If  $\pi_0$  is uniform, the density of  $\pi$  is proportional to  $L(\theta)$ , differing markedly from Jeffreys's prior density, instead proportional to  $\theta^{-1/2} (1 - \theta)^{-1/2}$  (Robert et al., 2012, p. 73), and the posterior density is proportional to the likelihood functions in Figure 1.

**Example 4.** Consider the two-state Markov chain  $X_{\theta,1}, X_{\theta,2}, \dots$  with probability transition matrix

$$\begin{bmatrix} 1 - \theta & \theta \\ \phi & 1 - \phi \end{bmatrix},$$

where  $\theta$  has a uniform blunted prior distribution  $\pi_0$  and  $\phi$  is known. The entropy rate of the process is

$$H(\theta) = \frac{\phi}{\theta + \phi} (-\log \Lambda(\theta)) + \frac{\theta}{\theta + \phi} (-\log \Lambda(\phi)),$$

where  $\Lambda(\bullet) = \bullet^\bullet (1 - \bullet)^{(1-\bullet)}$  (Cover and Thomas, 2006, pp. 73, 78), which is the whetted likelihood function

of Example 3. The whetted prior distribution is determined by

$$\pi(\theta) \propto L(\theta) = (\Lambda(\theta))^{\frac{\phi}{\theta+\phi}} (\Lambda(\phi))^{\frac{\theta}{\theta+\phi}}.$$

To extend the whetted distribution to continuous random variables, let  $Y_{\theta,1}, Y_{\theta,2}, \dots$  denote a stationary process such that, for all  $\theta \in \Theta$  and  $t = 1, 2, \dots$ , the distribution of  $Y_{\theta,t}$  is a Riemann-integrable probability density function  $g_{\theta,\tau}$  on an interval  $\mathcal{Y}$  for all  $\tau = 1, 2, \dots$ . The *differential whetted probability*  $\tilde{\pi}(\mathcal{H})$  that  $\vartheta \in \mathcal{H}$  is

$$\tilde{\pi}(\mathcal{H}) = \frac{\pi_0(\mathcal{H}) \int_{\mathcal{H}} d\pi_0(\theta|\mathcal{H}) 2^{-h(\theta)}}{\int d\pi_0(\theta) 2^{-h(\theta)}},$$

in which  $h(\theta)$  is the differential entropy rate  $h(\theta) = \lim_{\tau \rightarrow \infty} h(Y_{\theta,1}, \dots, Y_{\theta,\tau}) / \tau$ , where  $h(Y_{\theta,1}, \dots, Y_{\theta,\tau})$  is the differential entropy of  $Y_{\theta,1}, \dots, Y_{\theta,\tau}$  for every  $\theta \in \Theta$ . In analogy with  $L$ , the function  $\theta \mapsto \tilde{L}(\theta) = c2^{-h(\theta)}$  for any  $c > 0$  is called the *differential whetted likelihood function*.

**Example 5.** Suppose  $Y_{\theta,1}, Y_{\theta,2}, \dots$  are independent variates drawn from the normal distribution with unknown mean  $\mu$  and unknown standard deviation  $\sigma$ , where  $\theta = (\mu, \sigma)$ . Invariant measures that are limits of prior probability density functions include the left Haar measure and the right Haar measure, with densities proportional to  $\sigma^{-2}$  and  $\sigma^{-1}$ , respectively (Kass and Wasserman, 1996). Since  $2^{h(\mu,\sigma)} \propto \sigma$  (Michalowicz et al., 2013), the whetted likelihood is  $\tilde{L}(\mu, \sigma) \propto \sigma^{-1}$ . Thus, if  $\pi_0$  is the right Haar measure, then  $\tilde{\pi}$  is the left Haar measure, having a density proportional to  $\tilde{L}(\mu, \sigma) \sigma^{-1} \propto \sigma^{-2}$ . Using that  $\tilde{\pi}$  as the prior, the posterior predictive distribution of the next observation after observing  $n$  observations of sample mean  $\hat{\mu}$  and unbiased sample variance  $\hat{\sigma}^2$  is the  $t$  distribution with parameters  $\hat{\mu}$ ,  $(1 - n^{-1})\hat{\sigma}^2$ , and  $n - 1$  (Held and Sabanés Bové, 2014, p. 301). Since the posterior distribution has a mean of  $\hat{\mu}$ , the Bayes estimate of  $\mu$  under squared error loss is  $\hat{\mu}$  since that minimizes the posterior expected loss using  $\tilde{\pi}$  as the prior. Like  $\pi$  of Example 3, this  $\tilde{\pi}$  applies to small  $n$  as well as to large  $n$ .

The differential whetted probability is a limit of the following whetted probabilities that are themselves limits in the sense of Theorem 1.  $\mathcal{Y}_1, \mathcal{Y}_2, \dots$  is a sequence of interval subsets of  $\mathcal{Y}_\infty$  such that  $\mathcal{Y}_i \subset \mathcal{Y}_j$  for  $i < j$ , and  $\mathfrak{Y}_i(\Delta)$  is a partition of  $\mathcal{Y}_i$  into intervals of width  $\Delta > 0$  for  $i = 1, 2, \dots$  and for  $i = \infty$ , where



$\mathcal{Y}_\infty = \mathcal{Y}$ . Denote by  $X_{\theta,i,\Delta,t}$  a random member of  $\mathfrak{Y}_i(\Delta)$  with probability masses

$$P(X_{\theta,i,\Delta,t} = \mathcal{Y}') = P(X_{\theta,t} \in \mathcal{Y}' | X_{\theta,t} \in \mathcal{Y}_i)$$

for all  $\mathcal{Y}' \in \mathfrak{Y}_i(\Delta)$ ,  $\theta \in \Theta$ , and  $t = 1, 2, \dots$ . The entropy  $H(X_{\theta,i,\Delta,1}, \dots, X_{\theta,i,\Delta,\tau})$  and the whetted distribution for  $X_{\theta,i,\Delta,1}, X_{\theta,i,\Delta,2}, \dots$  are abbreviated by  $H(\theta, i, \Delta, \tau)$  and  $\pi_{i,\Delta}$ , respectively. That whetted distribution converges to  $\tilde{\pi}$ , the *differential whetted distribution*, in the following sense.

**Proposition 1.** *If  $\lim_{i \rightarrow \infty} H(\theta, i, \Delta, \tau) = H(\theta, \infty, \Delta, \tau)$  for all  $\Delta > 0$  and  $\tau = 1, 2, \dots$  and*

$$\lim_{\Delta \rightarrow 0} \lim_{i \rightarrow \infty} \lim_{\tau \rightarrow \infty} H(\theta, i, \Delta, \tau) + \log \Delta = \lim_{\tau \rightarrow \infty} \lim_{\Delta \rightarrow 0} \lim_{i \rightarrow \infty} H(\theta, i, \Delta, \tau) + \log \Delta$$

for all  $\theta \in \Theta$ , then, for all  $\mathcal{H} \in \mathfrak{H}$ ,

$$\lim_{\Delta \rightarrow 0} \lim_{i \rightarrow \infty} \pi_{i,\Delta}(\mathcal{H}) = \tilde{\pi}(\mathcal{H}).$$

*Proof.* Since  $g_{\theta,\tau}$  is Riemann integrable for all  $\theta \in \Theta$  and  $\tau = 1, 2, \dots$ ,

$$\lim_{\Delta \rightarrow 0} H(\theta, \infty, \Delta, \tau) + \log \Delta = h(Y_{\theta,1}, \dots, Y_{\theta,\tau})$$

(Cover and Thomas, 2006, Theorem 8.3.1). Thus,  $H(\theta, i, \Delta, \tau) + \log \Delta \rightarrow h(Y_{\theta,1}, \dots, Y_{\theta,\tau})$  and, with  $H(\theta, i, \Delta, \tau) / \tau \rightarrow H(\theta, i, \Delta)$ , which is the entropy rate of the process  $X_{\theta,i,\Delta,1}, X_{\theta,i,\Delta,2}, \dots$ ,

$$\begin{aligned} \lim_{\Delta \rightarrow 0} \lim_{i \rightarrow \infty} H(\theta, i, \Delta) + 0 &= \lim_{\Delta \rightarrow 0} \lim_{i \rightarrow \infty} \lim_{\tau \rightarrow \infty} H(\theta, i, \Delta, \tau) / \tau + \lim_{\Delta \rightarrow 0} \lim_{\tau \rightarrow \infty} (\log \Delta) / \tau \\ &= \lim_{\tau \rightarrow \infty} \left( \lim_{\Delta \rightarrow 0} \lim_{i \rightarrow \infty} H(\theta, i, \Delta, \tau) + \log \Delta \right) / \tau \\ &= \lim_{\tau \rightarrow \infty} h(Y_{\theta,1}, \dots, Y_{\theta,\tau}) / \tau = h(\theta). \end{aligned}$$

Substituting  $\lim_{\Delta \rightarrow 0} \lim_{i \rightarrow \infty} H(\theta, i, \Delta)$  for  $H(\theta)$  in equation (3) completes the proof.  $\square$

This research was partially supported by the Natural Sciences and Engineering Research Council of Canada (RGPIN/356018-2009) and by the Faculty of Medicine of the University of Ottawa.

## References

- Billingsley, P., 1999. *Convergence of Probability Measures*. Wiley Series in Probability and Statistics. Wiley, New York.
- Cover, T., Thomas, J., 2006. *Elements of Information Theory*. John Wiley & Sons, New York.
- Downey, R., 2008. Five lectures on algorithmic randomness. In: Chong, C., Feng, Q., Slaman, T., Woodin, W., Yang, Y. (Eds.), *Computational Prospects of Infinity: Tutorials*. World Scientific, Singapore, pp. 3–82.
- Downey, R., Hirschfeldt, D., 2010. *Algorithmic Randomness and Complexity. Theory and Applications of Computability*. Springer, New York.
- Edwards, A. W. F., 1992. *Likelihood*. Johns Hopkins Press, Baltimore.
- Held, L., Sabanés Bové, D., 2014. *Applied Statistical Inference* 10.
- Horibe, Y., 2003. A note on Kolmogorov complexity and entropy. *Applied Mathematics Letters* 16 (7), 1129–1130.
- Hutter, M., 2006. *Universal Artificial Intelligence: Sequential Decisions Based on Algorithmic Probability*. Texts in Theoretical Computer Science. An EATCS Series. Springer Berlin Heidelberg.
- Jeffreys, H., 1948. *Theory of Probability*. Oxford University Press, London.
- Kass, R. E., Wasserman, L., 1996. The selection of prior distributions by formal rules. *Journal of the American Statistical Association* 91, 1343–1370.
- Li, M., Vitányi, P., 2008. *An Introduction to Kolmogorov Complexity and Its Applications*. Texts in Computer Science. Springer New York.
- Michalowicz, J. V., Nichols, J. M., Bucholtz, F., 2013. *Handbook of Differential Entropy*. CRC Press, New York.
- Nies, A., 2012. *Computability and Randomness*. Oxford University Press, Oxford.
- Rathmanner, S., Hutter, M., 2011. A philosophical treatise of universal induction. *Entropy* 13 (6), 1076–1136.

Robert, C., Christensen, R., Johnson, W., Branscum, A., Hanson, T., 2012. Bayesian Ideas and Data Analysis.

Serfling, R. J., 1980. Approximation Theorems of Mathematical Statistics. Wiley, New York.

Solomonoff, R., 1978. Complexity-based induction systems: Comparisons and convergence theorems. IEEE Transactions on Information Theory 24 (4), 422–432.

Solomonoff, R. J., 2008. The probability of 'undefined' (non-converging) output in generating the universal probability distribution. Information Processing Letters 106 (6), 238–240.