# On the Kolmogorov Complexity of Binary Classifiers

Samuel Epstein\*

February 4, 2022

#### Abstract

We provide tight upper and lower bounds on the expected minimum Kolmogorov complexity of binary classifiers that are consistent with labeled samples. The expected size is not more than complexity of the target concept plus the conditional entropy of the labels given the sample.

## 1 Introduction

This paper provides bounds on the Kolmogorov complexity of binary classifiers. In machine learning, classification is the task of learning a binary function c from  $\mathbb{N}$  to bits  $\Sigma$ . The learner is given a sample consisting of pairs (x,b) for string x and bit b and outputs a binary classifier  $h: \mathbb{N} \to \Sigma$  that should match c as much as possible. Occam's razor says that "the simplest explanation is usually the best one." Simple hypothesis are resilient against overfitting to the sample data. In Section 2, we show some areas of machine learning which directly use the principle of Occam's razor. The question is, given a particular problem in machine learning, how simple can the hypotheses be?

In this paper, we provide tight upper and lower bounds of the length of the minimum description (a.k.a. Kolmogorov complexity,  $\mathbf{K}$ ) of hypotheses that are consistent with the labelled sample data.

We use a probabilistic model. The target concept is modeled by a random variable  $\mathcal{X}$  with distribution p over ordered lists of natural numbers. The random variable  $\mathcal{Y}$  models the labels, and has a distribution over lists of bits, where the distribution of  $\mathcal{X} \times \mathcal{Y}$  is p(x,y) with conditional probability requirement  $p(y|x) = \prod_{i=1..|x|} p(y_i|x_i)$ . Each such  $(x_i, y_i)$  is a labeled sample. A binary classifier f is consistent with labelled samples (x,y), if for all i,  $f(x_i) = y_i$ . Let  $\Gamma(x,y)$  be the minimum Kolmogorov complexity of a classifier consistent with (x,y).  $\mathcal{H}(\mathcal{Y}|\mathcal{X})$  is the conditional entropy of  $\mathcal{Y}$  given  $\mathcal{X}$ .

### Theorem.

- 1.  $\mathcal{H}(\mathcal{Y}|\mathcal{X}) \leq \mathbf{E}[\Gamma(\mathcal{X},\mathcal{Y})] <^{\log} \mathcal{H}(\mathcal{Y}|\mathcal{X}) + \mathbf{K}(p)$ .
- 2. For each  $c, b \in \mathbb{N}$ , there exists random labeled samples  $\mathcal{X} \times \mathcal{Y}$  with distribution p, such that, up to precision  $O(\log cb)$ ,  $\mathbf{E}[\Gamma(\mathcal{X}, \mathcal{Y})] = b + c$ ,  $\mathcal{H}(\mathcal{Y}|\mathcal{X}) = b$ , and  $\mathbf{K}(p) = c$ .

<sup>\*</sup>JP Theory Group. samepst@jptheorygroup.org

# 2 Related Work

There are many places the machine learning literature where Occam's razor is used. The goal of the extremely successful Minimum Description Length principle (see [GMP05]) is to find the most succinct hypothesis to model the data. This involves choosing a hypothesis H from a set of candidates that minimizes the length of the "code" for H, L(H), plus the length of the code describing the data D given H, denoted L(D|H). One difference between this paper and MDL is that MDL uses computable codes L, where we use the non-computable minimum program size.

The Occam Learning Algorithm [BEHW87] shows that a concept class is learnable if one can succinctly describe the training data. A concept c is a binary function over a finite set of strings. A concept class C is a set of such functions. The Occam Learning Algorithm takes in m labelled samples s of a concept c (a.k.a binary function) with encoding length  $\operatorname{Size}(c)$  and returns a hypothesis h consistent with c on s with  $\operatorname{Size}(h) \leq (n \times \operatorname{Size}(c))^{\alpha} m^{\beta}$ , for some  $\alpha \geq 0$  and  $0 \leq \beta < 1$ , where n is max length among the samples s. If a concept class has an Occam Algorithm, then it is PAC-learnable [Val84] and has a finite VC dimension [Vap98], which means that the concept class can be efficiently learned. Thus succinct representation of the hypothesis is connected to learnability of a target concept.

The structural risk minimization principle in statistical learning theory [Vap98] looks for the optimal relationship between the quality of the approximation of the target concept by a hypothesis and the complexity class which the hypothesis is in (a.k.a VC dimension). Thus simply described hypothesis will be more valued than complicated ones.

### Algorithmic Information Theory

In [Eps21], it was shown that the minimum description length of a binary classifier consistent with n samples is less than n plus the amount of information that the samples have with the halting sequence. Theorem 1 generalizes this result to clopen sets and computable measures.

The study of Kolmogorov complexity originated from the work of [Kol65]. The canonical self-delimiting form of Kolmogorov complexity was introduced in [ZL70] and treated later in [Cha75]. More information about the history of the concepts used in this paper can be found the text-book [LV08].

Theorem 1 in this paper is an inequality including the mutual information of the encoding of a finite set with the halting sequence. A history of the origin of the mutual information of a string with the halting sequence can be found in [VV04].

A string is stochastic if it is typical of a simple elementary probability distribution. A string is typical of a probability measure if it has a low deficiency of randomness. The notion of the deficiency of randomness with respect to a measure follows from the work of [She83], and also studied in [KU87, V'Y87, She99]. Aspects involving stochastic objects were studied in [She83, She99, V'Y87, V'Y99].

# 3 Conventions

Let  $\Sigma$ ,  $\Sigma^*$ , and  $\Sigma^{\infty}$  be the sets of bits, finite strings, and infinite strings. We use  $\langle x \rangle$  to represent a self delimiting code for  $x \in \Sigma^*$ , with  $\langle x \rangle = 1^{\|x\|} 0x$ . For sets  $C \subseteq \Sigma^{\infty}$  and  $D \subseteq \Sigma^*$ ,  $C \trianglelefteq D = \{x : x \in D, \Gamma_x \subseteq C\}$ . A clopen set in the Cantor space is a finite union og intervals. For clopen set C,  $\langle C \rangle = \langle \{x : \Gamma_x \text{ is maximal in } S\} \rangle$ . The indicator function of a mathematical statement A is denoted by [A], where if A is true then [A] = 1, otherwise [A] = 0.

For positive real functions f the terms  $<^+ f$ ,  $>^+ f$ , and  $=^+ f$  represent < f + O(1), > f - O(1), and  $= f \pm O(1)$ , respectively. For nonnegative real function f the terms  $<^{\log} f$ ,  $>^{\log} f$  and  $=^{\log} f$  represent  $< f + O(\log(f+1))$ ,  $> f - O(\log(f+1))$ , and  $= \pm O(\log(f+1))$ , respectively.

A probability measure Q over  $\mathbb{N}$  is elementary if  $|\{a:Q(a)>0\}|<\infty$  and Range(Q) consists of all rationals. Elementary measures Q can be encoded into finite strings  $\langle Q \rangle$ .

We use a universal prefix free algorithm U, where we say  $U_{\alpha}(x) = y$  if U, on main input x and auxiliary input  $\alpha$ , outputs y. We define Kolmogorov complexity with respect to U, with for  $x \in \Sigma^*$ ,  $y \in \Sigma^{*\infty}$ ,  $\mathbf{K}(x/y) = \min\{\|p\| : U_y(p) = x\}$ . By the chain rule,  $\mathbf{K}(x,y) = \mathbf{K}(x) + \mathbf{K}(y/x,\mathbf{K}(x))$ . The halting sequence  $\mathcal{H} \in \Sigma^{\infty}$  is the unique infinite sequence where  $\mathcal{H}[i] = [U(i) \text{ halts}]$ . The information that  $x \in \Sigma^*$  has about  $\mathcal{H}$ , conditional to  $y \in \Sigma^* \cup \Sigma^{\infty}$ , is  $\mathbf{I}(x;\mathcal{H}/y) = \mathbf{K}(x|y) - \mathbf{K}(x/\langle y,\mathcal{H} \rangle)$ . The Kolmogorov complexity of an infinite sequence  $\alpha \in \Sigma^{\infty}$  is the size of the smallest input to U which will output, without halting,  $\alpha$  on the output tape.

This paper uses notions of stochasticity in the field of algorithmic statistics [VS17]. A string x is stochastic, i.e. has a low  $\mathbf{Ks}(x)$  score if it is typical of a simple probability distribution. The deficiency of randomness function of a string x with respect to an elementary probability measure P is  $\mathbf{d}(x|P) = |-\log P(x)| - \mathbf{K}(x|\langle P \rangle)$ .

```
Definition 1 (Stochasticity). For x, y \in \Sigma^*, \mathbf{Ks}(x) = \min\{\mathbf{K}(P) + 3\log \max\{\mathbf{d}(x|P), 1\} : P \text{ is an elementary probability measure}\}.
```

# 4 Results

Each binary classifier can be represented as an infinite sequence in the natural way. The following theorem is a statement about measures and clopen sets. It may be of independent interest, as it generalizes Theorem 6 from [Eps21].

#### Theorem 1.

```
For clopen set C \subseteq \Sigma^{\infty}, computable measure S, \min_{\alpha \in C} \mathbf{K}(\alpha) < \log -\log S(C) + \mathbf{I}(C; \mathcal{H}) + O(\mathbf{K}(S)).
```

**Proof.** Let  $s = \lceil -\log S(C) \rceil$ . We remove consideration of the complexity terms of S and s in the proof because of the size of the error terms of the theorem. Let P be an elementary probability measure that realizes  $\mathbf{Ks}(\langle C \rangle)$ . Let n be the maximum length of members of finite sets encoded in the support of P. More formally,  $n = \max\{\|x\| : x \in W \subset \Sigma^*, \langle W \rangle \in \operatorname{Supp}(P)\}$ . The max term can be used because P is elementary, and thus has a finite support.

The randomness deficiency of S with respect to P is  $d = \max\{\mathbf{d}(C|P), 1\}$ . Let  $c \in \mathbb{N}$  be a constant solely dependent on the universal Turing machine U to be determined later. Let  $\kappa$  be a probability measure over lists L of  $cd2^s$  strings of length n, where  $\kappa(L) = \prod_{i=1}^{cd2^s} S(L_i)$ . Let  $\mathbf{i}(W, L)$ 

be an indicator function over sets  $W \subset \Sigma^*$  and lists  $L \subseteq \Sigma^n$ , with  $\mathbf{i}(W, L) = [S(W) \ge 2^{-s}, W \le L = \varnothing]$ .

$$\mathbf{E}_{L \sim \kappa} \mathbf{E}_{\langle W \rangle \sim P} [\mathbf{i}(W, L)] \leq \sum_{\text{clopen } W \subseteq \Sigma^{\infty}} P(\langle W \rangle) (1 - 2^{-s})^{cd2^{s}} \leq e^{-2^{-s}cd2^{s}} = e^{-cd}.$$

Thus there exists a list L of  $cd2^s$  strings such that  $\mathbf{E}_{\langle W \rangle \sim P}[\mathbf{i}(W,L)] < e^{-cd}$ . This L can be found with brute force search, with  $\mathbf{K}(L|c,d,P) = O(1)$ . Using L, we can define the following P-test,  $t(W) = e^{cd}\mathbf{i}(W,L)$ , with  $\sum_W P(W)t(W) \leq 1$ . It must be that  $C \leq L \neq \emptyset$ , otherwise t will give C a high score, with  $t(C) = e^{cd}$ . This causes the following contradiction for large enough c solely dependent on the universal Turing machine U, with

$$\mathbf{K}(C|c,d,\langle P\rangle) < -\log t_L(C)P(\langle C\rangle) + O(1)$$

$$\mathbf{K}(C|c,d,\langle P\rangle) < -\log P(\langle C\rangle) - (\lg e)cd + O(1)$$

$$(\lg e)cd < -\log P(\langle C\rangle) - \mathbf{K}(C|\langle P\rangle) + \mathbf{K}(d,c) + O(1)$$

$$(\lg e)cd < d + \mathbf{K}(d,c) + O(1).$$

We roll c into the additive constants of the rest of the proof. So there exists  $x \in C \subseteq L$ , with

$$\mathbf{K}(x) <^{+} \log |L| + \mathbf{K}(L)$$

$$<^{+} \log |L| + \mathbf{K}(d, P)$$

$$<^{+} \log d + s + \mathbf{K}(d) + \mathbf{K}(P)$$

$$<^{+} s + \mathbf{Ks}(\langle C \rangle).$$

Since  $x \in C \subseteq L$ ,  $\Gamma_x \subseteq C$ . Thus there is a program g that outputs x and then an infinite sequence of 0's. Since  $x0^{\infty} \in C$  and  $||g|| <^+ \mathbf{K}(x)$ ,

$$\min_{\alpha \in C} \mathbf{K}(\alpha) \le ||g|| <^+ \mathbf{K}(x) <^+ s + \mathbf{Ks}(\langle C \rangle).$$

Using Lemma 10 in [Eps21], which states  $\mathbf{Ks}(x) < \log \mathbf{I}(x; \mathcal{H})$ , we get the final form of the proof,

$$\min_{\alpha \in C} \mathbf{K}(\alpha) <^{\log} s + \mathbf{I}(C; \mathcal{H}) + O(\mathbf{K}(S)). \tag{1}$$

The following lemma is perhaps surprising because it shows that the  $\mathbf{I}(\cdot;\mathcal{H})$  terms in inequalities can be removed by averaging over a computable probability.

**Lemma 1.** For computable probability p,  $\sum_{x} p(x) \mathbf{I}(x; \mathcal{H}) <^{+} \mathbf{K}(p)$ .

**Proof.** This follows from Theorem 3.1.3 in [G21], and we will reproduce its arguments. Since  $\mathbf{K}(x/\mathcal{H})$  is the length of a self delimiting code,

$$\sum_{x} p(x) \mathbf{K}(x/\mathcal{H}) \ge \mathcal{H}(p),$$

where  $\mathcal{H}(p)$  is the entropy of p. Furthermore, for all  $x \in \Sigma^*$ ,  $\mathbf{K}(x) <^+ -\log p(x) + \mathbf{K}(p)$ . Therefore

$$\sum_{x} p(x)\mathbf{K}(x) <^{+} \sum_{x} p(x)(-\log p(x)) + \mathbf{K}(p) <^{+} \mathcal{H}(p) + \mathbf{K}(p).$$

So

$$\sum_{x} p(x)\mathbf{I}(x;\mathcal{H}) = \sum_{x} p(x) \left(\mathbf{K}(x) - \mathbf{K}(x/\mathcal{H})\right) <^{+} \mathcal{H}(p) + \mathbf{K}(p) - \sum_{x} p(x)\mathbf{K}(x/\mathcal{H}) <^{+} \mathbf{K}(p).$$

The following theorem is the main result of the theorem. It essentially involves averaging the inequality of Theorem 1 over the target probability. The  $I(\cdot; \mathcal{H})$  term vanishes due to Lemma 1. The theorem is slightly better than the simpler statement in the introduction. The terms in the theorem are defined in the introduction.

Theorem 2. 
$$\mathbf{E}[\Gamma(\mathcal{X}, \mathcal{Y})] < \mathcal{H}(\mathcal{Y}|\mathcal{X}) + \mathbf{K}(p) + O(\log \mathcal{H}(\mathcal{Y}|\mathcal{X})).$$

**Proof.** Binary classifiers are identified by infinite sequences  $\alpha \in \Sigma^{\infty}$ . We define the computable measure S, where  $S(x) = \prod_{n=1..|x|} p(x_n|n)$ , where  $\mathbf{K}(S|p) = O(1)$ . Let  $\{(x_i, y_i)\}$  be a set of labelled samples and we define clopen set  $C_{x,y} = \{\alpha : \alpha \in \Sigma^{\infty}, \alpha[x_i] = y_i\}$ . Then  $S(C_{x,y}) = p(y|x)$ . By Theorem 1, relativized to p,

$$\min_{\alpha \in C_{x,y}} \mathbf{K}(\alpha|p) <^{\log} - \log S(C_{x,y}) + \mathbf{I}(C_{x,y}; \mathcal{H}|p) + O(\mathbf{K}(S|p))$$

$$<^{\log} - \log S(C_{x,y}) + \mathbf{I}(C_{x,y}; \mathcal{H}|p)$$

$$<^{\log} - \log p(y|x) + \mathbf{I}(C_{x,y}; \mathcal{H}|p)$$

$$\sum_{x,y} p(x,y) \min_{\alpha \in C_{x,y}} \mathbf{K}(\alpha|p) <^{\log} \sum_{x,y} p(x,y)(-\log p(y|x)) + \sum_{x,y} p(x,y)\mathbf{I}(C_{x,y}; \mathcal{H}|p).$$
(3)

Applying Lemma 1 relative to p, we get

$$\sum_{x,y} p(x,y) \mathbf{I}(C_{x,y}; \mathcal{H}|p) <^{+} \mathbf{K}(p|p) = O(1).$$

$$\tag{4}$$

Combining equations 3 and 4,

$$\sum_{x,y} p(x,y) \min_{\alpha \in C_{x,y}} \mathbf{K}(\alpha|p) < \sum_{x,y} p(x,y) (-\log p(y|x)) + O\left(\sum_{x,y} p(x,y) \log(-\log p(y|x))\right)$$

$$< \sum_{x,y} p(x,y) (-\log p(y|x)) + O\left(\log \sum_{x,y} p(x,y) (-\log p(y|x))\right)$$

$$\mathbf{E}[\Gamma(\mathcal{X},\mathcal{Y})] - \mathbf{K}(p) < \mathcal{H}(\mathcal{Y}|\mathcal{X}) + O(\log \mathcal{H}(\mathcal{Y}|\mathcal{X}))$$

$$\mathbf{E}[\Gamma(\mathcal{X},\mathcal{Y})] < \mathcal{H}(\mathcal{Y}|\mathcal{X}) + \mathbf{K}(p) + O(\log \mathcal{H}(\mathcal{Y}|\mathcal{X})).$$

The following theorems provides lower bounds to the Kolmogorov complexity of classifiers.

Theorem 3.  $\mathcal{H}(\mathcal{Y}|\mathcal{X}) \leq \mathbf{E}[\Gamma(\mathcal{X},\mathcal{Y})]$ 

**Proof.**  $\mathbf{E}[\Gamma(\mathcal{X},\mathcal{Y})] = \sum_{x} p(x) \sum_{y} p(y|x) \Gamma(x,y)$ . For a fixed x, ranged over y,  $\Gamma(x,y)$  represents the length of a self-delimiting code. Due to properties of conditional entropy,  $\sum_{x} p(x) \sum_{y} p(y|x) \Gamma(x,y) \geq \sum_{x} p(x) \sum_{y} p(y|x) (-\log p(y|x)) = \mathcal{H}(\mathcal{Y}|\mathcal{X})$ .

**Theorem 4.** For each  $c, b \in \mathbb{N}$ , there exists random labeled samples  $\mathcal{X} \times \mathcal{Y}$  with distribution p, such that, up to precision  $O(\log cb)$ ,

- 1.  $\mathcal{H}(\mathcal{Y}|\mathcal{X}) = b$ ,
- 2.  $\mathbf{K}(p) = c$ ,
- 3.  $\mathbf{E}[\Gamma(\mathcal{X}, \mathcal{Y})] = b + c$ .

**Proof.** We ignore all  $O(\log cd)$  terms. So equality = is equivalent to =  $\pm O(\log cd)$ . We define a probability p(x,y) over the first n = 2c + 2b + 2 numbers and corresponding bits. Thus we can describe p as a probability measure over strings of size n, making sure to maintain p's conditional probability restriction described in the introduction.

Let  $z \in \Sigma^c$  be a random string of size c, with  $c <^+ \mathbf{K}(z)$ . For all strings  $w \in \Sigma^b$  of size b,  $p(\langle z \rangle \langle w \rangle) = 2^{-b}$ , with  $\|\langle z \rangle \langle w \rangle\| = n$ .  $\mathcal{H}(\mathcal{Y}|\mathcal{X}) = -\sum_{w \in \Sigma^b} 2^{-b} (\log p(\langle w \rangle \langle z \rangle)) = -\sum_{w \in \Sigma^b} 2^{-b} (\log 2^{-b}) = b$ . Furthermore  $\mathbf{K}(p) = c$ .

The infinite sequence  $\alpha = \langle w \rangle \langle z \rangle 0^{\infty}$  realizes  $\Gamma(\langle w \rangle \langle z \rangle)$  up to an additive constant for each  $w \in \Sigma^b$ . Thus  $\mathbf{K}(\alpha) = \mathbf{K}(z, w)$ .  $\mathbf{E}[\Gamma(\mathcal{X}, \mathcal{Y})] = 2^{-b} \sum_{w \in \Sigma^b} \mathbf{K}(\langle z \rangle \langle w \rangle) = \mathbf{K}(z) + 2^{-b} \sum_{w \in \Sigma^b} \mathbf{K}(w/z, \mathbf{K}(z))$ . Using Theorem 3.1.3 in [G21] conditioned on  $\langle z, \mathbf{K}(z) \rangle$ , we get that  $\sum_{w \in \Sigma^b} 2^{-b} \mathbf{K}(w/z, \mathbf{K}(z)) = \mathcal{H}(\mathcal{U}_b) \pm \mathbf{K}(b/z, \mathbf{K}(z)) = b$ , where  $\mathcal{U}_b$  is the uniform measure over strings of size b. So  $\mathbf{E}[\Gamma(\mathcal{X}, \mathcal{Y})] = \mathbf{K}(z) + b = b + c$ .

## 5 Discussion

The results of this paper are for binary classifiers that are completely consistent with the sample data. One area of research is looking into the description length of classifiers that have a small classification error. In general case, the bounds of Theorem 2 cannot be improved upon. However, we hypothesize that there exist interesting sets of models p whose expected classifier description length is much smaller than  $\mathbf{K}(p)$ . Another open research area is the intersection of algorithmic information theory with other areas of statistical learning theory, including density estimation and regression.

# References

- [BEHW87] A. Blumer, A. Ehrenfeucht, D. Haussler, and Warmuth. Occam's razor. *Information processing letters*, 24(6):377–380, 1987.
- [Cha75] G. J. Chaitin. A Theory of Program Size Formally Identical to Information Theory. Journal of the ACM, 22(3):329–340, 1975.
- [Eps21] Samuel Epstein. All sampling methods produce outliers. *IEEE Transactions on Information Theory*, 67(11):7568–7578, 2021.
- [G21] Peter Gács. Lecture notes on descriptional complexity and randomness. CoRR, abs/2105.04704, 2021.

- [GMP05] P. Grunwald, I. Myung, and M. Pitt. Advances in minimum description length theory and applications. MIT Press, Cambridge, MA, USA, 2005.
- [Kol65] A. N. Kolmogorov. Three approaches to the quantitative definition of information. Problems in Information Transmission, 1:1–7, 1965.
- [KU87] A. N. Kolmogorov and V. A. Uspensky. Algorithms and Randomness. SIAM Theory of Probability and Its Applications, 32(3):389–412, 1987.
- [LV08] M. Li and P. Vitányi. An Introduction to Kolmogorov Complexity and Its Applications. Springer Publishing Company, Incorporated, 3 edition, 2008.
- [She83] A. Shen. The concept of (alpha,beta)-stochasticity in the Kolmogorov sense, and its properties. *Soviet Mathematics Doklady*, 28(1):295–299, 1983.
- [She99] A. Shen. Discussion on Kolmogorov Complexity and Statistical Analysis. *The Computer Journal*, 42(4):340–342, 1999.
- [Val84] L. Valiant. A theory of the learnable. Commun. ACM, 27(11):1134–1142, 1984.
- [Vap98] V. Vapnik. Statistical Learning Theory. Wiley-Interscience, Hoboken, NJ, 1998.
- [VS17] Nikolay K. Vereshchagin and Alexander Shen. Algorithmic statistics: Forty years later. In *Computability and Complexity*, pages 669–737, 2017.
- [VV04] N. Vereshchagin and P. Vitányi. Kolmogorov's Structure Functions and Model Selection. *IEEE Transactions on Information Theory*, 50(12):3265 3290, 2004.
- [V'Y87] V.V. V'Yugin. On Randomness Defect of a Finite Object Relative to Measures with Given Complexity Bounds. SIAM Theory of Probability and Its Applications, 32:558–563, 1987.
- [V'Y99] V.V. V'Yugin. Algorithmic complexity and stochastic properties of finite binary sequences. *The Computer Journal*, 42:294–317, 1999.
- [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 Math. Surveys*, page 11, 1970.