

Universal probability-free conformal prediction

Vladimir Vovk and Dusko Pavlovic

March 20, 2016

Abstract

We construct a universal prediction system in the spirit of Popper’s falsifiability and Kolmogorov complexity. This prediction system does not depend on any statistical assumptions, but under the IID assumption it dominates, although in a rather weak sense, conformal prediction.

Not for nothing do we call the laws of nature “laws”:
the more they prohibit, the more they say.

The Logic of Scientific Discovery
KARL POPPER

1 Introduction

In this paper we consider the problem of predicting labels, assumed to be binary, of a sequence of objects. This is an online version of the standard problem of binary classification. Namely, we will be interested in infinite sequences of observations

$$\omega = (z_1, z_2, \dots) = ((x_1, y_1), (x_2, y_2), \dots) \in (\mathbf{X} \times 2)^\infty$$

(also called *infinite data sequences*), where \mathbf{X} is an *object space* and $2 := \{0, 1\}$. For simplicity, we will assume that \mathbf{X} is a given finite set of, say, binary strings (the intuition being that finite objects can always be encoded as binary strings).

Finite sequences $\sigma \in (\mathbf{X} \times 2)^*$ of observations will be called *finite data sequences*. If σ_1, σ_2 are two finite data sequences, their concatenation will be denoted (σ_1, σ_2) ; σ_2 is also allowed to be an element of $\mathbf{X} \times 2$. A standard partial order on $(\mathbf{X} \times 2)^*$ is defined as follows: $\sigma_1 \sqsubseteq \sigma_2$ means that σ_1 is a prefix of σ_2 ; $\sigma_1 \sqsubset \sigma_2$ means that $\sigma_1 \sqsubseteq \sigma_2$ and $\sigma_1 \neq \sigma_2$.

We use the notation $\mathbb{N} := \{1, 2, \dots\}$ for the set of positive integers and $\mathbb{N}_0 := \{0, 1, 2, \dots\}$ for the set of nonnegative integers. If $\omega \in (\mathbf{X} \times 2)^\infty$ and $n \in \mathbb{N}_0$, $\omega^n \in (\mathbf{X} \times 2)^n$ is the prefix of ω of length n .

A *situation* is a concatenation $(\sigma, x) \in (\mathbf{X} \times 2)^* \times \mathbf{X}$ of a finite data sequence σ and an object x ; our task in the situation (σ, x) is to be able to predict the

label of the new object x given the sequence σ of labelled objects. Given a situation $s = (\sigma, x)$ and a label $y \in 2$, we let (s, y) stand for the finite data sequence $(\sigma, (x, y))$, which is the concatenation of s and y .

2 Laws of nature as prediction systems

According to Popper’s [1] view of the philosophy of science, scientific laws of nature should be falsifiable: if a finite sequence of observations contradicts such a law, we should be able to detect it. (Popper often preferred to talk about scientific theories or statements instead of laws of nature.) The empirical content of a law of nature is the set of its potential falsifiers ([1], Sections 31 and 35). We start from formalizing this notion in our toy setting, interpreting the requirement that we should be able to detect falsification as that we should be able to detect it eventually.

Formally, we define a *law of nature* L to be a recursively enumerable prefix-free subset of $(\mathbf{X} \times 2)^*$ (where *prefix-free* means that $\sigma_2 \notin L$ whenever $\sigma_1 \in L$ and $\sigma_1 \sqsubset \sigma_2$). Intuitively, these are the potential falsifiers, i.e., sequences of observations prohibited by the law of nature. The requirement of being recursively enumerable is implicit in the notion of a falsifier, and the requirement of being prefix-free reflects the fact that extensions of prohibited sequences of observations are automatically prohibited and there is no need to mention them in the definition.

A law of nature L gives rise to a prediction system: in a situation $s = (\sigma, x)$ it predicts that the label $y \in 2$ of the new object x will be an element of

$$\Pi_L(s) := \{y \in 2 \mid (s, y) \notin L\}. \quad (1)$$

There are three possibilities in each situation s :

- The law of nature makes a prediction, either 0 or 1, in situation s when the prediction set (1) is of size 1, $|\Pi_L(s)| = 1$.
- The prediction set is empty, $|\Pi_L(s)| = 0$, which means that the law of nature has been falsified.
- The law of nature refrains from making a prediction when $|\Pi_L(s)| = 2$. This can happen in two cases:
 - the law of nature was falsified in past: $\sigma' \in L$ for some $\sigma' \sqsubseteq \sigma$;
 - the law of nature has not been falsified as yet.

3 Strong prediction systems

The notion of a law of nature is static; experience tells us that laws of nature eventually fail and are replaced by other laws. Popper represented his picture of this process by formulas (“evolutionary schemas”) similar to

$$\text{PS}_1 \rightarrow \text{TT}_1 \rightarrow \text{EE}_1 \rightarrow \text{PS}_2 \rightarrow \dots \quad (2)$$

(introduced in his 1965 talk on which [2], Chapter 6, is based and also discussed in several other places in [2] and [3]; in our notation we follow Wikipedia). In response to a problem situation PS, a tentative theory TT is subjected to attempts at error elimination EE, whose success leads to a new problem situation PS and scientists come up with a new tentative theory TT, etc. In our toy version of this process, tentative theories are laws of nature, problem situations are situations in which our current law of nature becomes falsified, and there are no active attempts at error elimination (so that error elimination simply consists in waiting until the current law of nature becomes falsified).

If L and L' are laws of nature, we define $L \sqsubset L'$ to mean that for any $\sigma' \in L'$ there exists $\sigma \in L$ such that $\sigma \sqsubset \sigma'$. To formalize the philosophical picture (2), we define a *strong prediction system* \mathcal{L} to be a nested sequence $L_1 \sqsubset L_2 \sqsubset \dots$ of laws of nature L_1, L_2, \dots that are jointly recursively enumerable, in the sense of the set $\{(\sigma, n) \in (\mathbf{X} \times 2)^* \times \mathbb{N} \mid \sigma \in L_n\}$ being recursively enumerable.

The interpretation of a strong prediction system $\mathcal{L} = (L_1, L_2, \dots)$ is that L_1 is the initial law of nature used for predicting the labels of new objects until it is falsified; as soon as it is falsified we start looking for and then using for prediction the following law of nature L_2 until it is falsified in its turn, etc. Therefore, the prediction set in a situation $s = (\sigma, x)$ is natural to define as the set

$$\Pi_{\mathcal{L}}(s) := \{y \in 2 \mid (s, y) \notin \cup_{n=1}^{\infty} L_n\}. \quad (3)$$

As before, it is possible that $\Pi_{\mathcal{L}}(s) = \emptyset$.

Fix a situation $s = (\sigma, x) \in (\mathbf{X} \times 2)^* \times \mathbf{X}$. Let $n = n(s)$ be the largest integer such that s has a prefix in L_n . It is possible that $n = 0$ (when s does not have such prefixes), but if $n \geq 1$, s will also have prefixes in L_{n-1}, \dots, L_1 , by the definition of a strong prediction system. Then L_{n+1} will be the current law of nature; all earlier laws, L_n, L_{n-1}, \dots, L_1 , have been falsified. The prediction (3) in situation s is then interpreted as the set of all observations y that are not prohibited by the current law L_{n+1} .

In the spirit of the theory of Kolmogorov complexity, we would like to have a universal prediction system. However, we are not aware of any useful notion of a universal strong prediction system. Therefore, in the next section we will introduce a wider notion of a prediction system that does not have this disadvantage.

4 Weak prediction systems and universal prediction

A *weak prediction system* \mathcal{L} is defined to be a sequence (not required to be nested in any sense) L_1, L_2, \dots of laws of nature $L_n \subseteq (\mathbf{X} \times 2)^*$ that are jointly recursively enumerable.

Remark. Popper's evolutionary schema (2) was the simplest one that he con-

sidered; his more complicated ones, such as

$$\begin{array}{c}
 \nearrow \text{TT}_a \rightarrow \text{EE}_a \rightarrow \text{PS}_{2a} \rightarrow \dots \\
 \text{PS}_1 \rightarrow \text{TT}_b \rightarrow \text{EE}_b \rightarrow \text{PS}_{2b} \rightarrow \dots \\
 \searrow \text{TT}_c \rightarrow \text{EE}_c \rightarrow \text{PS}_{2c} \rightarrow \dots
 \end{array}$$

(cf. [2], pp. 243 and 287), give rise to weak rather than strong prediction systems.

In the rest of this paper we will omit “weak” in “weak prediction system”. The most basic way of using a prediction system \mathcal{L} for making a prediction in situation $s = (\sigma, x)$ is as follows. Decide on the maximum number N of errors you are willing to make. Ignore all L_n apart from L_1, \dots, L_N in \mathcal{L} , so that the prediction set in situation s is

$$\Pi_{\mathcal{L}}^N(s) := \{y \in 2 \mid \forall n \in \{1, \dots, N\} : (s, y) \notin L_n\}.$$

Notice that this way we are guaranteed to make at most N mistakes: making a mistake eliminates at least one law in the list $\{L_1, \dots, L_N\}$.

Similarly to the usual theory of conformal prediction, another way of packaging \mathcal{L} 's prediction in situation s is, instead of choosing the threshold (or *level*) N in advance, to allow the user to apply her own threshold: in a situation s , for each $y \in 2$ report the attained level

$$\pi_{\mathcal{L}}^s(y) := \min \{n \in \mathbb{N} \mid (s, y) \in L_n\} \quad (4)$$

(with $\min \emptyset := \infty$). The user whose threshold is N will then consider $y \in 2$ with $\pi_{\mathcal{L}}^s(y) \leq N$ as prohibited in s . Notice that the function (4) is upper semicomputable (for a fixed \mathcal{L}).

The strength of a prediction system $\mathcal{L} = (L_1, L_2, \dots)$ at level N is determined by its *N-part*

$$\mathcal{L}_{\leq N} := \bigcup_{n=1}^N L_n.$$

At level N , the prediction system L prohibits $y \in 2$ as continuation of a situation s if and only if $(s, y) \in \mathcal{L}_{\leq N}$.

The following lemma says that there exists a universal prediction system, in the sense that it is stronger than any other prediction system if we ignore a multiplicative increase in the number of errors made.

Lemma 1. *There is a universal prediction system \mathcal{U} , in the sense that for any prediction system \mathcal{L} there exists a constant $C > 0$ such that, for any N ,*

$$\mathcal{L}_{\leq N} \subseteq \mathcal{U}_{\leq CN}. \quad (5)$$

Proof. Let $\mathcal{L}^1, \mathcal{L}^2, \dots$ be a recursive enumeration of all prediction systems; their component laws of nature will be denoted $(L_1^k, L_2^k, \dots) := \mathcal{L}^k$. For each $n \in \mathbb{N}$,

define the n th component U_n of $\mathcal{U} = (U_1, U_2, \dots)$ as follows. Let the binary representation of n be

$$(a, 0, 1, \dots, 1), \quad (6)$$

where a is a binary string (starting from 1) and the number of 1s in the $1, \dots, 1$ is $k - 1 \in \mathbb{N}_0$ (this sentence is the definition of $a = a(n)$ and $k = k(n)$ in terms of n). If the binary representation of n does not contain any 0s, a and k are undefined, and we set $U_n := \emptyset$. Otherwise, set

$$U_n := L_A^k,$$

where $A \in \mathbb{N}$ is the number whose binary representation is a . In other words, \mathcal{U} consists of the components of \mathcal{L}^k , $k \in \mathbb{N}$; namely, L_1^k is placed in \mathcal{U} as $U_{3 \times 2^{k-1} - 1}$ and then L_2^k, L_3^k, \dots are placed at intervals of 2^k :

$$U_{3 \times 2^{k-1} - 1 + 2^k(i-1)} = L_i^k, \quad i = 1, 2, \dots$$

It is easy to see that

$$\mathcal{L}_{\leq N}^k \subseteq \mathcal{U}_{\leq 3 \times 2^{k-1} - 1 + 2^k(N-1)}, \quad (7)$$

which is stronger than (5). \square

Let us fix a universal prediction system \mathcal{U} . By $K(\mathcal{L})$ we will denote the smallest prefix complexity of the programs for computing a prediction system \mathcal{L} . The following lemma makes (5) uniform in \mathcal{L} showing how C depends on \mathcal{L} .

Lemma 2. *There is a constant $C > 0$ such that, for any prediction system \mathcal{L} and any N , the universal prediction system \mathcal{U} satisfies*

$$\mathcal{L}_{\leq N} \subseteq \mathcal{U}_{\leq C 2^{K(\mathcal{L})} N}. \quad (8)$$

Proof. Follow the proof of Lemma 1 replacing the “code” $(0, 1, \dots, 1)$ for \mathcal{L}^k in (6) by any prefix-free description of \mathcal{L}^k (with its bits written in the reverse order). Then the modification

$$\mathcal{L}_{\leq N}^k \subseteq \mathcal{U}_{\leq 2^{k'+1} - 1 + 2^{k'}(N-1)}$$

of (7) with $k' := K(\mathcal{L}^k)$ implies that (8) holds for some universal prediction system, which, when combined with the statement of Lemma 1, implies that (8) holds for our chosen universal prediction system \mathcal{U} . \square

This is a corollary for laws of nature:

Corollary 1. *There is a constant C such that, for any law of nature L , the universal prediction system \mathcal{U} satisfies*

$$L \subseteq \mathcal{U}_{\leq C 2^{K(L)}}. \quad (9)$$

Proof. We can regard laws of nature L to be a special case of prediction systems identifying L with $\mathcal{L} := (L, L, \dots)$. It remains to apply Lemma 2 to \mathcal{L} setting $N := 1$. \square

We can equivalently rewrite (5), (8), and (9) as

$$\Pi_{\mathcal{U}}^{C_N}(s) \subseteq \Pi_{\mathcal{L}}^N(s), \quad (10)$$

$$\Pi_{\mathcal{U}}^{C^{2^{K(\mathcal{L})}}N}(s) \subseteq \Pi_{\mathcal{L}}^N(s), \quad (11)$$

and

$$\Pi_{\mathcal{U}}^{C^{2^{K(L)}}}(s) \subseteq \Pi_L(s), \quad (12)$$

respectively, for all situations s . Intuitively, (10) says that the prediction sets output by the universal prediction system are at least as precise as the prediction sets output by any other prediction system \mathcal{L} if we ignore a constant factor in specifying the level N ; and (11) and (12) indicate the dependence of the constant factor on \mathcal{L} .

5 Universal conformal prediction under the IID assumption

Comparison of prediction systems and conformal predictors is hampered by the fact that the latter are designed for the case where we have a constant amount of noise for each observation, and so we expect the number of errors to grow linearly rather than staying bounded. In this situation a reasonable prediction set is $\Pi_{\mathcal{L}}^{\epsilon N}(s)$, where N is the number of observations in the situation s . For a small ϵ using $\Pi_{\mathcal{L}}^{\epsilon N}(s)$ means that we trust the prediction system whose percentage of errors so far is at most ϵ .

Up to this point our exposition has been completely probability-free, but in the rest of this section we will consider the special case where the data are generated in the IID manner. For simplicity, we will only consider computable conformity measures that take values in the set \mathbb{Q} of rational numbers.

Corollary 2. *Let Γ be a conformal predictor based on a computable conformity measure taking values in \mathbb{Q} . Then there exists $C > 0$ such that, for almost all infinite sequences of observations $\omega = ((x_1, y_1), (x_2, y_2), \dots) \in (\mathbf{X} \times \mathbf{Y})^\infty$ and all significance levels $\epsilon \in (0, 1)$, from some N on we will have*

$$\Pi_{\mathcal{U}}^{C_N \epsilon \ln^2(1+1/\epsilon)}((\omega^N, x_{N+1})) \subseteq \Gamma^\epsilon((\omega^N, x_{N+1})). \quad (13)$$

This corollary asserts that the prediction set output by the universal prediction system is at least as precise as the prediction set output by Γ if we increase slightly the significance level: from ϵ to $C\epsilon \ln^2(1+1/\epsilon)$. It involves not just multiplying by a constant (as is the case for (5) and (8)–(12)) but also the logarithmic term $\ln^2(1+1/\epsilon)$.

It is easy to see that we can replace the C in (13) by $C^{2^{K(\Gamma)}}$, where C now does not depend on Γ (and $K(\Gamma)$ is the smallest prefix complexity of the programs for computing the conformity measure on which Γ is based).

Proof of Corollary 2. Let

$$\epsilon' := 2^{\lceil \log \epsilon \rceil + 1},$$

where \log stands for the base 2 logarithm. (Intuitively, we simplify ϵ , in the sense of Kolmogorov complexity, by replacing it by a number of the form 2^{-m} for an integer m , and make it at least twice as large as the original ϵ .) Define a prediction system (both weak and strong) \mathcal{L} as, essentially, $\Gamma^{\epsilon'}$; formally, $\mathcal{L} := (L_1, L_2, \dots)$ and L_n is defined to be the set of all ω^N , where ω ranges over the infinite data sequences and N over \mathbb{N} , such that the set

$$\left\{ i \in \{1, \dots, N\} \mid y_i \notin \Gamma^{\epsilon'}((\omega^{i-1}, x_i)) \right\}$$

is of size n and contains N . The prediction system \mathcal{L} is determined by ϵ' , so that $K(\mathcal{L})$ does not exceed (apart from the usual additive constant) $K(\epsilon')$. By the standard validity property of conformal predictors ([6], Corollary 1.1), Hoeffding's inequality, and the Borel–Cantelli lemma,

$$\Pi_{\mathcal{L}}^{\epsilon'^N}((\omega^N, x_{N+1})) \subseteq \Gamma^{\epsilon}((\omega^N, x_{N+1})) \quad (14)$$

from some N on almost surely. By Lemma 2 (in the form of (11)),

$$\Pi_{\mathcal{U}}^{C_1 2^{K(\epsilon')} \epsilon'^N}((\omega^N, x_{N+1})) \subseteq \Pi_{\mathcal{L}}^{\epsilon'^N}((\omega^N, x_{N+1})) \quad (15)$$

for all N . The statement (13) of the corollary is obtained by combining (14), (15), and

$$2^{K(\epsilon')} \leq C_2 \ln^2(1 + 1/\epsilon).$$

To check the last inequality, remember that $\epsilon' = 2^{-m}$ for an integer m , which we assume to be positive, without loss of generality; therefore, our task reduces to checking that

$$2^{K(m)} \leq C_3 \ln^2(1 + 2^m),$$

i.e.,

$$2^{K(m)} \leq C_4 m^2.$$

Since $2^{-K(m)}$ is the universal semimeasure on the positive integers (see, e.g., [5], Theorem 7.29), we even have

$$2^{K(m)} \leq C_5 m(\log m)(\log \log m) \cdots (\log \cdots \log m),$$

where the product contains all factors that are greater than 1 (see [4], Appendix A). \square

6 Conclusion

In this note we have ignored the computational resources, first of all, the required computation time and space (memory). Developing versions of our definitions and results taking into account the time of computations is a natural next step. In analogy with the theory of Kolmogorov complexity, we expect that the simplest and most elegant results will be obtained for computational models that are more flexible than Turing machines, such as Kolmogorov–Uspensky algorithms and Schönhage machines.

Acknowledgments.

We thank the anonymous referees for helpful comments. This work has been supported by the Air Force Office of Scientific Research (grant “Semantic Completions”), EPSRC (grant EP/K033344/1), and the EU Horizon 2020 Research and Innovation programme (grant 671555).

References

- [1] Karl R. Popper. *Logik der Forschung*. Springer, Vienna, 1934. English translation: *The Logic of Scientific Discovery*. Hutchinson, London, 1959.
- [2] Karl R. Popper. *Objective Knowledge: An Evolutionary Approach*. Clarendon Press, Oxford, revised edition, 1979. First edition: 1972.
- [3] Karl R. Popper. *All Life is Problem Solving*. Abingdon, Routledge, 1999.
- [4] Jorma Rissanen. A universal prior for integers and estimation by minimum description length. *Annals of Statistics*, 11:416–431, 1983.
- [5] Alexander Shen. Around Kolmogorov complexity: Basic notions and results. In Vladimir Vovk, Harris Papadopoulos, and Alexander Gammerman, editors, *Measures of Complexity: Festschrift for Alexey Chervonenkis*, chapter 7, pages 75–115. Springer, Cham, 2015.
- [6] Vladimir Vovk. The basic conformal prediction framework. In Vineeth N. Balasubramanian, Shen-Shyang Ho, and Vladimir Vovk, editors, *Conformal Prediction for Reliable Machine Learning: Theory, Adaptations, and Applications*, chapter 1, pages 3–19. Elsevier, Amsterdam, 2014.