

# Universal Well-Calibrated Algorithm for On-line Classification

Vladimir Vovk



практические выводы  
теории вероятностей  
могут быть обоснованы  
в качестве следствий  
гипотез о *предельной*  
при данных ограничениях  
сложности изучаемых явлений

**On-line Compression Modelling Project**

Working Paper #3

November 7, 2002

Project web site:  
<http://vovk.net/kp>

# Abstract

We study the problem of on-line classification in which the prediction algorithm is given a “confidence level”  $1 - \delta$  and is required to output as its prediction a range of labels (intuitively, those labels deemed compatible with the available data at the level  $\delta$ ) rather than just one label; as usual, the examples are assumed to be generated independently from the same probability distribution  $P$ . The prediction algorithm is said to be “well-calibrated” for  $P$  and  $\delta$  if the long-run relative frequency of errors does not exceed  $\delta$  almost surely w.r. to  $P$ . For well-calibrated algorithms we take the number of “uncertain” predictions (i.e., those containing more than one label) as the principal measure of predictive performance. The main result of this paper is the construction of a prediction algorithm which, for any (unknown)  $P$  and any  $\delta$ : (a) makes errors independently and with probability  $\delta$  at every trial (in particular, is well-calibrated for  $P$  and  $\delta$ ); (b) makes in the long run no more uncertain predictions than any other prediction algorithm that is well-calibrated for  $P$  and  $\delta$ ; (c) processes example  $n$  in time  $O(\log n)$ .

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Main result</b>	<b>3</b>
<b>3</b>	<b>Construction of a universal well-calibrated region predictor</b>	<b>6</b>
3.1	Preliminaries . . . . .	6
3.2	Transductive Confidence Machines . . . . .	7
3.3	Universal TCM . . . . .	7
<b>4</b>	<b>Fine details of region prediction</b>	<b>9</b>
<b>5</b>	<b>Conclusion</b>	<b>11</b>
<b>A</b>	<b>Appendix: Proofs</b>	<b>13</b>
A.1	Proof sketch of Proposition 2 . . . . .	13
A.2	Proof sketch of Proposition 4 . . . . .	15
A.3	Proof sketch of Proposition 5 . . . . .	18

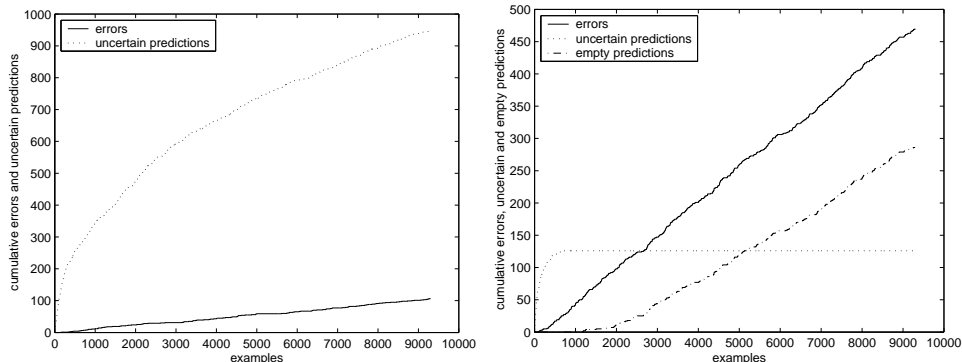


Figure 1: TCM at 99% (left) and 95% (right) on the USPS data set

## 1 Introduction

Typical machine learning algorithms output a point prediction for the label of an unknown object. This paper continues study of an algorithm, called Transductive Confidence Machine (TCM) and introduced in [11, 7], that complements its predictions with some measures of confidence. There are different ways of presenting TCM’s output; in this paper (as in the related [10, 9]) we use TCM as a “region predictor”, in the sense that for any confidence level  $1 - \delta$  it outputs a predictive region rather than a point prediction.

Paper [10] shows that any TCM is well-calibrated when used in the on-line mode: for any confidence level  $1 - \delta$  the long-run relative frequency of erroneous predictions does not exceed  $\delta$ . What makes this property of TCM especially appealing is that it is far from being just an asymptotic phenomenon: a slight modification of TCM called randomized TCM (rTCM; randomization is needed to break ties and deal efficiently with borderline cases) makes errors independently at different trials and with probability  $\delta$  at each trial; well-calibratedness then immediately follows by the Borel strong law of large numbers.

The justification of the study of TCM given in [10] was its good performance on real-world and standard benchmark data sets. For example, Figure 1 shows that for the standard confidence levels 99% and 95% most examples in the well-known USPS data set (randomly permuted) can be predicted categorically (by a simple Nearest Neighbors TCM): the predictive region contains only one label.

This paper presents theoretical results about TCM’s performance; we

show that there exists a *universal* rTCM, which, for any confidence level  $1 - \delta$  and without knowing the true distribution  $P$  generating the examples:

- produces, asymptotically, no more uncertain predictions than any other prediction algorithm that is well-calibrated for  $P$  and  $\delta$ ;
- produces, asymptotically, at least as many empty predictions as any other prediction algorithm that is well-calibrated for  $P$  and  $\delta$  and whose percentage of uncertain predictions is optimal (in the sense of the previous item).

The importance of the first item is obvious: we want to minimize the number of uncertain predictions. This criterion ceases to work, however, when the number of uncertain predictions stabilizes, as in Figure 1 (right). In such cases the number of empty predictions becomes important: empty predictions (automatically leading to an error) provide a warning that the object is untypical (looks very different from the previous objects), and one would like to be warned as often as possible, taking into account that the relative frequency of errors (including empty predictions) is guaranteed not to exceed  $\delta$  in the long run.

This paper’s result elaborates on [9], where it was shown that an optimal TCM exists when the distribution  $P$  generating the examples is known. Here we consider only randomized TCM, so we drop the adjective “randomized”.

The two areas of mainstream machine learning that are most closely connected with this paper are PAC learning theory and Bayesian learning theory. Whereas we often use the rich arsenal of mathematical tools developed in these fields, they do not provide the same kind of guarantees (a prespecified probability of error, with errors at different trials independent) under unknown  $P$ ; for more details, see [10] and references therein. Several papers (such as [6, 5]) extend the standard PAC framework by allowing the prediction algorithm to abstain from making a prediction at some trials. Our results show that for any confidence level  $1 - \delta$  there exists a prediction algorithm that: (a) makes a wrong prediction with relative frequency at most  $\delta$ ; (b) has an optimal frequency of abstentions among the prediction algorithms that satisfy property (a) (for details, see Remark 2 on p. 5). Paper [5] is especially close to the approach of this paper, defining a very natural TCM in the situation where a hypothesis class is given (the “empirical log ratio” of [5], taken with appropriate sign, can be used as “individual strangeness measure”, as defined in §3).

## 2 Main result

In our learning protocol, Nature outputs pairs  $(x_1, y_1), (x_2, y_2), \dots$  called *examples*. Each example  $(x_i, y_i)$  consists of an *object*  $x_i$  and its *label*  $y_i$ ; the objects are chosen from a measurable space  $\mathbf{X}$  called the *object space* and the labels are elements of a measurable space  $\mathbf{Y}$  called the *label space*. In this paper we assume that  $\mathbf{Y}$  is finite (and endowed with the  $\sigma$ -algebra of all subsets). The protocol includes variables  $\text{Err}_n$  (the total number of errors made up to and including trial  $n$ ) and  $\text{err}_n$  (the binary variable showing whether an error is made at trial  $n$ ); it also includes analogous variables  $\text{Unc}_n, \text{unc}_n, \text{Emp}_n, \text{emp}_n$  for uncertain and empty predictions:

```

Err0 := 0; Unc0 := 0; Emp0 := 0;
FOR n = 1, 2, ...:
  Nature outputs xn ∈ X;
  Predictor outputs Γn ⊆ Y;
  Nature outputs yn ∈ Y;
  errn := { 1 if yn ∉ Γn ; Errn := Errn-1 + errn;
            0 otherwise
  uncn := { 1 if |Γn| > 1 ; Uncn := Uncn-1 + uncn;
            0 otherwise
  empn := { 1 if |Γn| = 0 ; Empn := Empn-1 + empn
            0 otherwise
END FOR.

```

We will use the notation  $\mathbf{Z} := \mathbf{X} \times \mathbf{Y}$  for the *example space*;  $\Gamma_n$  will be called the *predictive region* (or just *prediction*).

We will be assuming that each example  $z_n = (x_n, y_n)$ ,  $n = 1, 2, \dots$ , is output according to a probability distribution  $P$  in  $\mathbf{Z}$  and the examples are independent of each other (so the sequence  $z_1 z_2 \dots$  is output by the power distribution  $P^\infty$ ). This is Nature's randomized strategy.

A *region predictor* is a measurable function

$$\Gamma_\gamma(x_1, \tau_1, y_1, \dots, x_{n-1}, \tau_{n-1}, y_{n-1}, x_n, \tau_n),$$

where  $n = 1, 2, \dots$ , the  $(x_i, y_i) \in \mathbf{Z}$ ,  $i = 1, \dots, n-1$ , are examples,  $x_n \in \mathbf{X}$  is an object, and  $\tau_i \in [0, 1]$  ( $i = 1, \dots, n$ ), which satisfies

$$\begin{aligned} & \Gamma_{\gamma_1}(x_1, \tau_1, y_1, \dots, x_{n-1}, \tau_{n-1}, y_{n-1}, x_n, \tau_n) \\ & \subseteq \Gamma_{\gamma_2}(x_1, \tau_1, y_1, \dots, x_{n-1}, \tau_{n-1}, y_{n-1}, x_n, \tau_n) \end{aligned}$$

whenever  $\gamma_1 \leq \gamma_2$ . Since we are interested in prediction with confidence, the region predictor is given an extra input  $\gamma = 1 - \delta \in (0, 1)$ , which we call the *confidence level* (typically it is close to 1, standard values being 99% and 95%); the complementary value  $\delta$  is called the *significance level*. We will always assume that  $\tau_n$  are independent random variables uniformly distributed in  $[0, 1]$ ; this makes a region predictor a family (indexed by  $\gamma \in (0, 1)$ ) of Predictor's randomized strategies.

We will often use the notation  $\text{err}_n$ ,  $\text{unc}_n$ , etc., in the case where Predictor and Nature are using given randomized strategies: for example,  $\text{err}_n(P^\infty, \Gamma_{1-\delta})$  is the random variable equal to 1 if Predictor is right at trial  $n$  and 0 otherwise. It is always assumed that the random numbers  $\tau_n$  used by  $\Gamma$  and the random examples  $z_n$  chosen by Nature are independent.

We say that a region predictor  $\Gamma$  is (conservatively) *well-calibrated* for a probability distribution  $P$  in  $\mathbf{Z}^\infty$  and a significance level  $\delta \in (0, 1)$  if

$$\limsup_{n \rightarrow \infty} \frac{\text{Err}_n(P^\infty, \Gamma_{1-\delta})}{n} \leq \delta \quad \text{a.s.}$$

We say (following [9]) that  $\Gamma$  is *optimal* for  $P$  and  $\delta$  if, for any region predictor  $\Gamma^\dagger$  which is well-calibrated for  $P$  and  $\delta$ ,

$$\limsup_{n \rightarrow \infty} \frac{\text{Unc}_n(P^\infty, \Gamma_{1-\delta})}{n} \leq \liminf_{n \rightarrow \infty} \frac{\text{Unc}_n(P^\infty, \Gamma_{1-\delta}^\dagger)}{n} \quad \text{a.s.} \quad (1)$$

(the symbol ‘‘a.s.’’ in such a context always assumes that the random numbers used by  $\Gamma$  and  $\Gamma^\dagger$  are independent). Of course, the definition of optimality is natural only for well-calibrated  $\Gamma$ .

A region predictor  $\Gamma$  is *universal well-calibrated* if:

- it is well-calibrated for any  $P$  and  $\delta$ ;
- it is optimal for any  $P$  and  $\delta$ ;
- for any  $P$ , any  $\delta$ , and any region predictor  $\Gamma^\dagger$  which is well-calibrated and optimal for  $P$  and  $\delta$ ,

$$\liminf_{n \rightarrow \infty} \frac{\text{Emp}_n(P^\infty, \Gamma_{1-\delta})}{n} \geq \limsup_{n \rightarrow \infty} \frac{\text{Emp}_n(P^\infty, \Gamma_{1-\delta}^\dagger)}{n} \quad \text{a.s.}$$

Recall that a measurable space  $\mathbf{X}$  is *Borel* if it is isomorphic to a measurable subset of the interval  $[0, 1]$ . The class of Borel spaces is very rich; for example, all Polish spaces (such as finite-dimensional Euclidean spaces  $\mathbb{R}^n$ ,  $\mathbb{R}^\infty$ , functional spaces  $C$  and  $D$ ) are Borel.

**Theorem 1** *Suppose the object space  $\mathbf{X}$  is Borel. There exists a universal well-calibrated region predictor.*

This is the main result of the paper; in §3 we construct a universal well-calibrated region predictor (processing example  $n$  in time  $O(\log n)$ ) and in §4 outline the idea of the proof that it indeed satisfies the required properties.

**Remark** In this paper we are interested in the theoretical properties of region predictors for a fixed confidence level. This does not mean, however, that fixing a confidence level in advance is the right thing to do in practice; at the very least, two or more conventional levels should be used. For example, we could say that the prediction is “highly certain” if  $|\Gamma_{0.99}| \leq 1$  and “certain” if  $|\Gamma_{0.95}| \leq 1$ ; similarly, we could say that the new object (whose label is being predicted) is “highly untypical” if  $|\Gamma_{0.99}| = 0$  and “untypical” if  $|\Gamma_{0.95}| = 0$ . To avoid the dependence on arbitrarily chosen conventional levels, the range of possible predictive regions  $\Gamma_{1-\delta}$ ,  $\delta \in (0, 1)$ , can be summarized by reporting the *confidence*

$$\sup\{\gamma : |\Gamma_\gamma| \leq 1\},$$

the *credibility*

$$\inf\{\delta : |\Gamma_{1-\delta}| = 0\},$$

and the *prediction*  $\Gamma_\gamma$ , where  $\gamma$  is the confidence ( $\Gamma_\gamma$  is non-empty for TCM and usually contains exactly one label). Reporting the prediction, confidence, and credibility, as in [11, 7], is analogous to reporting the observed level of significance ([2], p. 66) in statistics.

**Remark** The protocol of [6, 5] is in fact a restriction of our protocol, in which Predictor is only allowed to output a one-element set or the whole of  $\mathbf{Y}$ ; the latter is interpreted as abstention. (And in the situation where  $\text{Err}_n$  and  $\text{Unc}_n$  are of primary interest, as in this paper, the difference between these two protocols is not very significant.) The universal well-calibrated region predictor can be adapted to the restricted protocol by replacing an uncertain prediction with  $\mathbf{Y}$  and replacing an empty prediction with a randomly chosen label. In this way we obtain a prediction algorithm in the restricted protocol which is well-calibrated and has an optimal frequency of abstentions, in the sense of (1), among the well-calibrated algorithms.

### 3 Construction of a universal well-calibrated region predictor

#### 3.1 Preliminaries

If  $\tau$  is a number in  $[0, 1]$ , we split it into two numbers  $\tau', \tau'' \in [0, 1]$  as follows: if the binary expansion of  $\tau$  is  $0.a_1a_2\dots$  (redefine the binary expansion of 1 to be  $0.11\dots$ ), set  $\tau' := 0.a_1a_3a_5\dots$  and  $\tau'' := 0.a_2a_4a_6\dots$ . If  $\tau$  is distributed uniformly in  $[0, 1]$ , then both  $\tau'$  and  $\tau''$  are, and they are independent of each other.

We will often apply our procedures (e.g., the “individual strangeness measure” in §3.2, the Nearest Neighbors rule in §3.3) not to the original objects  $x \in \mathbf{X}$  but to *extended objects*  $(x, \sigma) \in \tilde{\mathbf{X}} := \mathbf{X} \times [0, 1]$ , where  $x$  is complemented by a random number  $\sigma$  (to be extracted from one of the  $\tau_n$ ). In other words, along with examples  $(x, y)$  we will also consider *extended examples*  $(x, \sigma, y) \in \mathbf{Z} := \mathbf{X} \times [0, 1] \times \mathbf{Y}$ .

Let us set  $\mathbf{X} := [0, 1]$ ; we can do this without loss of generality since  $\mathbf{X}$  is Borel. This makes the extended object space  $\tilde{\mathbf{X}} = [0, 1]^2$  a linearly ordered set with the lexicographic order:  $(x_1, \sigma_1) < (x_2, \sigma_2)$  means that either  $x_1 = x_2$  and  $\sigma_1 < \sigma_2$  or  $x_1 < x_2$ . We say that  $(x_1, \sigma_1)$  is *nearer* to  $(x_3, \sigma_3)$  than  $(x_2, \sigma_2)$  is if

$$|x_1 - x_3, \sigma_1 - \sigma_3| < |x_2 - x_3, \sigma_2 - \sigma_3|, \quad (2)$$

where

$$|x, \sigma| := \begin{cases} (x, \sigma) & \text{if } (x, \sigma) \geq (0, 0) \\ (-x, -\sigma) & \text{otherwise.} \end{cases}$$

Our construction will be based on the Nearest Neighbors algorithm, which is known to be strongly universally consistent in the traditional theory of pattern recognition (see, e.g., [4], Chapter 11); the random components  $\sigma$  are needed for tie-breaking. As usual, to give a precise meaning to the expression “the  $k$ th nearest neighbor” in a sequence of extended objects to another extended object  $v$ , we use the convention that extended objects with smaller indices in the sequence are considered to be nearer to  $v$  (this particular convention is not essential, since adding the random components  $\sigma$  ensures that ties will occur with probability zero).



## 3.2 Transductive Confidence Machines

TCM is a way of transition from what we call an “individual strangeness measure” to a region predictor. A family of measurable functions  $\{A_n : n = 1, 2, \dots\}$ , where  $A_n : \mathbf{Z}^n \rightarrow \mathbb{R}^n$  for all  $n$ , is called an *individual strangeness measure* if, for any  $n = 1, 2, \dots$ , each  $\alpha_i$  in

$$A_n : (w_1, \dots, w_n) \mapsto (\alpha_1, \dots, \alpha_n) \quad (3)$$

is determined by  $w_i$  and the multiset  $\wr w_1, \dots, w_n \wr$ .

The *TCM associated with an individual strangeness measure*  $A_n$  is the following region predictor  $\Gamma_{1-\delta}(x_1, \tau_1, y_1, \dots, x_{n-1}, \tau_{n-1}, y_{n-1}, x_n, \tau_n)$ : at any trial  $n$  and for any label  $y \in \mathbf{Y}$ , define

$$(\alpha_1, \dots, \alpha_n) := A_n((x_1, \tau'_1, y_1), \dots, (x_{n-1}, \tau'_{n-1}, y_{n-1}), (x_n, \tau'_n, y)),$$

and include  $y$  in  $\Gamma_{1-\delta}$  if and only if

$$\tau_n'' < \frac{\#\{i = 1, \dots, n : \alpha_i \geq \alpha_n\} - n\delta}{\#\{i = 1, \dots, n : \alpha_i = \alpha_n\}} \quad (4)$$

(in particular, include  $y$  in  $\Gamma_{1-\delta}$  if  $\#\{i = 1, \dots, n : \alpha_i > \alpha_n\}/n > \delta$  and do not include  $y$  in  $\Gamma_{1-\delta}$  if  $\#\{i = 1, \dots, n : \alpha_i \geq \alpha_n\}/n \leq \delta$ ).

A *TCM* is the TCM associated with some individual strangeness measure. It was shown in [10] that

**Proposition 1** *Every TCM is well-calibrated for every  $P$  and  $\delta$ .*

## 3.3 Universal TCM

Fix a monotonically non-decreasing sequence of integer numbers  $K_n$ ,  $n = 1, 2, \dots$ , such that

$$K_n \rightarrow \infty, K_n = o(n/\ln n) \quad (5)$$

as  $n \rightarrow \infty$ . The *Nearest Neighbors TCM* is defined as follows. Let  $w_1, \dots, w_n$  be a sequence of extended examples  $w_i = (x_i, \sigma_i, y_i)$ . To define the corresponding  $\alpha$ s (see (3)), we first define Nearest Neighbors approximations  $P_n^\neq(y | x_i, \sigma_i)$  to the true (but unknown) conditional probabilities  $P(y | x_i)$ : for every extended example  $(x_i, \sigma_i, y_i)$  in the sequence,

$$P_n^\neq(y | x_i, \sigma_i) := N^\neq(x_i, \sigma_i, y)/K_n, \quad (6)$$

where  $N^\neq(x_i, \sigma_i, y)$  is the number of  $j = 1, \dots, n$  such that  $y_j = y$  and  $(x_j, \sigma_j)$  is one of the  $K_n$  nearest neighbors of  $(x_i, \sigma_i)$  in the sequence  $((x_1, \sigma_1), \dots, (x_{i-1}, \sigma_{i-1}), (x_{i+1}, \sigma_{i+1}), \dots, (x_n, \sigma_n))$ . (The upper index  $\neq$  reminds us of the fact that  $(x_i, \sigma_i)$  is not counted as one of its own nearest neighbors in this definition.) If  $K_n \geq n$  or  $K_n \leq 0$ , this definition does not work, so set, e.g.,  $P_n^\neq(y | x_i, \sigma_i) := 1/|\mathbf{Y}|$  for all  $y$  and  $i$  (this particular convention is not essential since, by (5),  $0 < K_n < n$  from some  $n$  on).

Define the “empirical predictability function”  $f_n^\neq$  by

$$f_n^\neq(x_i, \sigma_i) := \max_{y \in \mathbf{Y}} P_n^\neq(y | x_i, \sigma_i).$$

For each  $(x_i, \sigma_i)$  fix some

$$\hat{y}_n(x_i, \sigma_i) \in \arg \max_y P_n^\neq(y | x_i, \sigma_i)$$

(e.g., take the first element of  $\arg \max_y P_n^\neq(y | x_i, \sigma_i)$  in a fixed ordering of  $\mathbf{Y}$ ) and define the mapping (3) (where  $w_i = (x_i, \sigma_i, y_i)$ ,  $i = 1, \dots, n$ ) setting

$$\alpha_i := \begin{cases} -f_n^\neq(x_i, \sigma_i) & \text{if } y_i = \hat{y}_n(x_i, \sigma_i) \\ f_n^\neq(x_i, \sigma_i) & \text{otherwise.} \end{cases} \quad (7)$$

This completes the definition of the Nearest Neighbors TCM, which will later be shown to be universal.

**Proposition 2** *If  $\mathbf{X} = [0, 1]$  and  $K_n \rightarrow \infty$  sufficiently slowly, the Nearest Neighbors TCM can be implemented so that computations at trial  $n$  are performed in time  $O(\log n)$ .*

Proposition 2 assumes a computational model that allows operations (such as comparison) with real numbers. If  $\mathbf{X}$  is an arbitrary Borel space, for this proposition to be applicable  $\mathbf{X}$  should be imbedded in  $[0, 1]$  first; e.g., if  $\mathbf{X} \subseteq [0, 1]^n$ , an  $x = (x_1, \dots, x_n) \in \mathbf{X}$  can be represented as

$$(x_{1,1}, x_{2,1}, \dots, x_{n,1}, x_{1,2}, x_{2,2}, \dots, x_{n,2}, \dots) \in [0, 1],$$

where  $0.x_{i,1}x_{i,2}\dots$  is the binary expansion of  $x_i$ .

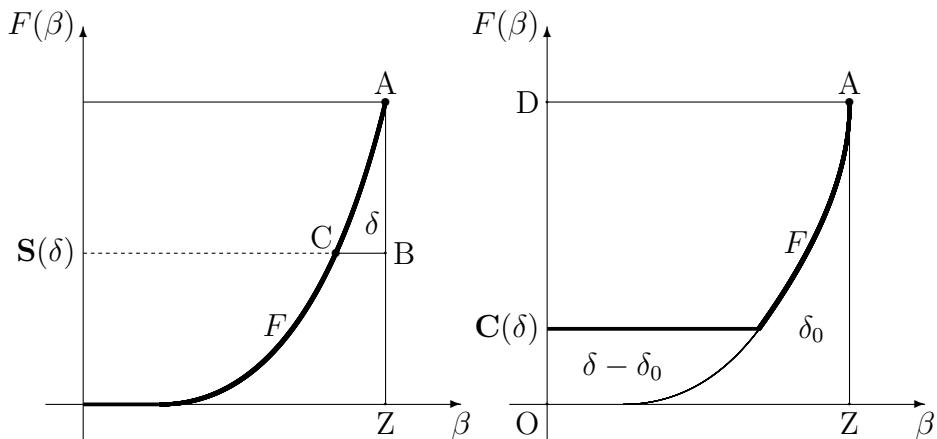


Figure 2: The predictability distribution function  $F$  and the success curve  $\mathbf{S}(\delta)$  (left); the complementary success curve  $\mathbf{C}(\delta)$  (right)

## 4 Fine details of region prediction

In this section we make first steps towards the proof of Theorem 1. Let  $P$  be the true distribution in  $\mathbf{Z}$  generating the examples. We denote by  $P_{\mathbf{X}}$  the marginal distribution of  $P$  in  $\mathbf{X}$  (i.e.,  $P_{\mathbf{X}}(E) := P(E \times \mathbf{Y})$ ) and by  $P_{\mathbf{Y}|\mathbf{X}}(y|x)$  the conditional probability that, for a random example  $(X, Y)$  chosen from  $P$ ,  $Y = y$  provided  $X = x$  (we fix arbitrarily a regular version of this conditional probability). We will often omit lower indices  $\mathbf{x}$  and  $\mathbf{y}|\mathbf{x}$  and  $P$  itself from our notation.

The *predictability* of an object  $x \in \mathbf{X}$  is

$$f(x) := \max_{y \in \mathbf{Y}} P(y|x)$$

and the *predictability distribution function* is the function  $F : [0, 1] \rightarrow [0, 1]$  defined by

$$F(\beta) := P\{x : f(x) \leq \beta\}.$$

An example of such a function  $F$  is given in Figure 2 (left), where the graph of  $F$  is the thick line.

The *success curve*  $\mathbf{S}$  of  $P$  is defined by the equality

$$\mathbf{S}_P(\delta) = \inf \left\{ B \in [0, 1] : \int_0^1 (F(\beta) - B)^+ d\beta \leq \delta \right\},$$

where  $t^+$  stands for  $\max(t, 0)$ ; the function  $\mathbf{S}$  is also of the type  $[0, 1] \rightarrow [0, 1]$ . Geometrically,  $\mathbf{S}_P(\delta)$  is defined from the graph of  $F$  as follows (see Figure 2, left; we usually drop the lower index  $P$ ): move the point B from A to Z until the area of the curvilinear triangle ABC becomes  $\delta$  or B reaches Z; the ordinate of B is then  $\mathbf{S}(\delta)$ .

The *complementary success curve*  $\mathbf{C}$  of  $P$  is defined by

$$\mathbf{C}_P(\delta) = \sup \left\{ B \in [0, 1] : B + \int_0^1 (F(\beta) - B)^+ d\beta \leq \delta \right\}.$$

Similarly to the case of  $\mathbf{S}(\delta)$ ,  $\mathbf{C}(\delta)$  is defined as the value such that the area of the part of the box AZOD below the thick line in Figure 2 (right) is  $\delta$  ( $\mathbf{C}(\delta) = 0$  if such a value does not exist).

Define the *critical significance level*  $\delta_0$  as

$$\delta_0 := \int_0^1 F(\beta) d\beta.$$

It is clear that

$$\begin{aligned} \delta \leq \delta_0 &\implies \int_0^1 (F(\beta) - \mathbf{S}(\delta))^+ d\beta = \delta \text{ \& } \mathbf{C}(\delta) = 0 \\ \delta \geq \delta_0 &\implies \mathbf{S}(\delta) = 0 \text{ \& } \mathbf{C}(\delta) + \int_0^1 (F(\beta) - \mathbf{C}(\delta))^+ d\beta = \delta. \end{aligned}$$

The following result is proven in [9].

**Proposition 3** *Let  $P$  be a probability distribution in  $\mathbf{Z}$  with success curve  $\mathbf{S}$  and  $\delta > 0$  be a significance level. If a region predictor  $\Gamma$  is well-calibrated for  $P$  and  $\delta$ , then*

$$\liminf_{n \rightarrow \infty} \frac{\text{Unc}_n(P^\infty, \Gamma_{1-\delta})}{n} \geq \mathbf{S}(\delta) \quad a.s.$$

In this paper we complement Proposition 3 with

**Proposition 4** *Let  $P$  be a probability distribution in  $\mathbf{Z}$  with success curve  $\mathbf{S}$  and complementary success curve  $\mathbf{C}$  and  $\delta > 0$  be a significance level. If a region predictor  $\Gamma$  is well-calibrated for  $P$  and  $\delta$  and satisfies*

$$\limsup_{n \rightarrow \infty} \frac{\text{Unc}_n(P^\infty, \Gamma_{1-\delta})}{n} \leq \mathbf{S}(\delta) \quad a.s., \quad (8)$$

then

$$\limsup_{n \rightarrow \infty} \frac{\text{Emp}_n(P^\infty, \Gamma_{1-\delta})}{n} \leq \mathbf{C}(\delta) \quad a.s.$$

Theorem 1 immediately follows from Propositions 1, 3, 4 and the following proposition.

**Proposition 5** *Suppose  $\mathbf{X}$  is Borel. The Nearest Neighbors TCM constructed in §3.3 satisfies, for any  $P$  and any significance level  $\delta$ ,*

$$\limsup_{n \rightarrow \infty} \frac{\text{Unc}_n(P^\infty, \Gamma_{1-\delta})}{n} \leq \mathbf{S}_P(\delta) \quad a.s. \quad (9)$$

and

$$\liminf_{n \rightarrow \infty} \frac{\text{Emp}_n(P^\infty, \Gamma_{1-\delta})}{n} \geq \mathbf{C}_P(\delta) \quad a.s. \quad (10)$$

## 5 Conclusion

We have shown that there exist universal well-calibrated region predictors, thus satisfying, to some degree, the desiderata mentioned in §1: well-calibratedness and universality. Notice, however, that the ways in which these two desiderata are satisfied are very different: the well-calibratedness holds in a very specific finitary sense, since the errors have probability  $\delta$  and are independent, whereas the universality is an asymptotic property.

## Acknowledgments

I am grateful to anonymous referees [STOC conference]. This work was partially supported by EPSRC (grant GR/R46670/01) BBSRC (grant 111/BIO14428) and EU (grant IST-1999-10226).

## References

- [1] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*. MIT Press, Cambridge, MA, second edition, 2001.
- [2] David R. Cox and David V. Hinkley. *Theoretical Statistics*. Chapman and Hall, London, 1974.
- [3] Luc Devroye, László Györfi, Adam Krzyżak, and Gábor Lugosi. On the strong universal consistency of nearest neighbor regression function estimates. *Annals of Statistics*, 22:1371–1385, 1994.

- [4] Luc Devroye, László Györfi, and Gábor Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer, New York, 1996.
- [5] Yoav Freund, Yishay Mansour, and Robert E. Schapire. Generalization bounds for averaged classifiers. July 2002. To appear in *Annals of Statistics*.
- [6] Ronald L. Rivest and R. Sloan. Learning complicated concepts reliably and usefully. In *Proceedings of the First Annual Conference on Computational Learning Theory*, pages 69–79, San Mateo, CA, 1988. Morgan Kaufmann.
- [7] Craig Saunders, Alex Gammerman, and Vladimir Vovk. Transduction with confidence and credibility. In *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence*, pages 722–726, 1999.
- [8] Albert N. Shiryaev. *Probability*. Springer, New York, second edition, 1996.
- [9] Vladimir Vovk. Asymptotic optimality of Transductive Confidence Machine, On-line Compression Modelling project, <http://vovk.net/kp>, Working Paper #2. In *Proceedings of the Thirteenth International Conference on Algorithmic Learning Theory*, volume 2533 of *Lecture Notes in Artificial Intelligence*, pages 336–350, 2002.
- [10] Vladimir Vovk. On-line Confidence Machines are well-calibrated, On-line Compression Modelling project, <http://vovk.net/kp>, Working Paper #1. In *Proceedings of the Forty Third Annual Symposium on Foundations of Computer Science*, pages 187–196. IEEE Computer Society, 2002.
- [11] Vladimir Vovk, Alex Gammerman, and Craig Saunders. Machine-learning applications of algorithmic randomness. In *Proceedings of the Sixteenth International Conference on Machine Learning*, pages 444–453, San Francisco, CA, 1999. Morgan Kaufmann.

# A Appendix: Proofs

## A.1 Proof sketch of Proposition 2

For simplicity we will be assuming that all extended objects  $(x_i, \tau'_i) \in [0, 1]^2$  are different (in any case, this is true with probability one, since  $\tau'_i$  are independent random numbers distributed as  $\mathbf{U}$ ). Our computational model has an operation of splitting  $\tau \in [0, 1]$  into  $\tau'$  and  $\tau''$  (or is allowed to generate both  $\tau'_n$  and  $\tau''_n$  at every trial  $n$ ).

We will use two main data structures in our implementation of the Nearest Neighbors TCM:

- a red-black binary *search tree* (see, e.g., [1], Chapters 13–15; the only two operations on red-black trees we need in this paper are the query SEARCH and the modifying operation INSERT);
- a growing *array* of nonnegative integer numbers indexed by numbers  $k \in \{-K_n, -K_n + 1, \dots, K_n\}$  (where  $n$  is the ordinal number of the example being processed).

Immediately after processing the  $n$ th extended example  $(x_n, \tau_n, y_n)$  the contents of these data structures are as follows:

- The search tree contains  $n$  vertices, corresponding to the extended examples  $(x_i, \tau_i, y_i)$  seen so far. The key of vertex  $i$  is the extended object  $(x_i, \tau'_i) \in [0, 1]^2$ ; the linear order on the keys is the lexicographic order. The other information contained in vertex  $i$  is the random number  $\tau''_i$ , the label  $y_i$ , the set  $\{P_n^\neq(y | x_i, \tau'_i) : y \in \mathbf{Y}\}$  of conditional probability estimates (6), the pointer to the following vertex (i.e., the vertex that has the smallest key greater than  $(x_i, \tau'_i)$ ; if there is no greater key, the pointer is NIL), and the pointer to the previous vertex (i.e., the vertex that has the greatest key smaller than  $(x_i, \tau'_i)$ ; if  $(x_i, \tau'_i)$  is the smallest key, the pointer is NIL).
- The array contains the numbers

$$N(k) := \# \{i = 1, \dots, n : \alpha_i = k/K_n\}.$$

Notice that the information contained in vertex  $i$  of the search tree is sufficient to find  $\hat{y}_n(x_i, \tau'_i)$  and  $\alpha_i$  in time  $O(1)$ .

We will say that an extended object  $(x_j, \tau'_j)$  is in the *vicinity* of an extended object  $(x_i, \tau'_i)$  if there are less than  $K_n$  extended objects  $(x_k, \tau'_k)$  (strictly) between  $(x_i, \tau'_i)$  and  $(x_j, \tau'_j)$ .

When a new object  $x_n$  becomes known, the algorithm does the following:

- Generates  $\tau'_n$  and  $\tau''_n$ .
- Locates the successor and predecessor of  $(x_n, \tau'_n)$  in the search tree (using the query SEARCH and the pointers to the following and previous vertices); this requires time  $O(\log n)$ .
- Computes the conditional probabilities  $\{P_n^\neq(y | x_n, \tau'_n) : y \in \mathbf{Y}\}$ ; this also gives  $\hat{y}_n(x_n, \tau'_n)$ . The required time is  $O(K_n) = O(\log n)$ .
- For each  $y \in \mathbf{Y}$  looks at what happens if the  $n$ th example is  $(x_n, \tau_n, y_n) = (x_n, \tau_n, y)$ : computes  $\alpha_n$  and updates (if necessary)  $\alpha_i$  for  $(x_i, \tau'_i)$  in the vicinity of  $(x_n, \tau'_n)$ ; using the array and  $\tau''_n$ , finds if  $y \in \Gamma_n$ . This requires time  $O(K_n^2) = O(\log n)$ .
- Outputs the predictive region  $\Gamma_n$  (time  $O(1)$ ).

When the label  $y_n$  arrives, the algorithm:

- Inserts the new vertex  $(x_n, \tau'_n, \tau''_n, y_n, \{P_n^\neq(y | x_n, \tau_n) : y \in \mathbf{Y}\})$  in the search tree, repairs the pointers to the following and previous elements for  $(x_n, \tau'_n)$ 's left and right neighbors, initializes the pointers to the following and previous elements for  $(x_n, \tau'_n)$  itself, and rebalances the tree (time  $O(\log n)$ ).
- Updates (if necessary) the conditional probabilities

$$\{P_{n-1}^\neq(y | x_i, \tau'_i) : y \in \mathbf{Y}\} \mapsto \{P_n^\neq(y | x_i, \tau'_i) : y \in \mathbf{Y}\}$$

for the  $2K_n$  existing vertices  $(x_i, \tau'_i)$  in the vicinity of  $(x_n, \tau'_n)$ ; this requires time  $O(K_n^2) = O(\log n)$ . The conditional probabilities for other  $(x_i, \tau'_i)$ ,  $i = 1, \dots, n-1$ , do not change.

- Updates the array, changing  $N(K_n \alpha_i)$  for the  $(x_i, \tau'_i) \neq (x_n, \tau'_n)$  in the vicinity of  $(x_n, \tau'_n)$  and for both old and new values of  $\alpha_i$  and changing  $N(K_n \alpha_n)$  (time  $O(K_n) = O(\log n)$ ).



In conclusion we discuss how to do updates required when  $K_n$  increases. An *epoch* is defined to be a maximal sequence of  $n$ s with the same  $K_n$ . Since the changes that need to be done when a new epoch starts are substantial, they will be spread over the whole preceding epoch. An epoch is *odd* if the corresponding  $K_n$  is odd and *even* if  $K_n$  is even. At every step in an epoch we prepare the ground for the next epoch. By the end of epoch  $n = A + 1, A + 2, \dots, B$  we need to change  $B$  sets  $\{P_n^\neq(y | x_i, \tau'_i) : y \in \mathbf{Y}\}$  in  $B - A$  steps (the duration of the epoch). Therefore, each vertex of the search tree should contain not only  $\{P_n^\neq(y | x_i, \tau'_i)\}$  for the current epoch but also  $\{P_n^\neq(y | x_i, \tau'_i)\}$  for the next epoch (two structures for holding  $\{P_n^\neq(y | x_i, \tau'_i)\}$  will suffice, one for even epochs and one for odd epochs). If  $K_n$  grows slowly enough (say, as  $\log n$ ),  $B/A = O(1)$ . At each step,  $O(1)$  sets  $\{P_n^\neq(y | x_i, \tau'_i)\}$  for the next epoch are added. This will take time  $O(K_n) = O(\log n)$ . As soon as a set  $\{P_n^\neq(y | x_i, \tau'_i)\}$  for the next epoch is added at some trial, both sets (for the current and next epoch) will have to be updated for each new example. In a similar way the array for the next epoch is gradually built up.

## A.2 Proof sketch of Proposition 4

The proof of Proposition 4 is similar to (but more complicated than) the proof of Theorems 1 and 1r in [9]; this proof sketch can be made rigorous using the Neyman-Pearson lemma, as in [9].

We will use the notations  $g'_{\text{left}}$  and  $g'_{\text{right}}$  for the left and right derivatives, respectively, of a function  $g$ . The following lemma parallels Lemma 2 in [9], which deals with  $\mathbf{S}(\delta)$ .

**Lemma 1** *The complementary success curve  $\mathbf{C} : [0, 1] \rightarrow [0, 1]$  always satisfies these properties:*

1. *There is a point  $\delta_0 \in [0, 1]$  (viz. the critical significance level) such that  $\mathbf{C}(\delta) = 0$  for  $\delta \leq \delta_0$  and  $\mathbf{C}(\delta)$  is concave for  $\delta \geq \delta_0$ .*
2.  *$\mathbf{C}'_{\text{right}}(\delta_0) < \infty$  and  $\mathbf{C}'_{\text{left}}(1) \geq 1$ ; therefore, for  $\delta \in (\delta_0, 1)$ ,  $1 \leq \mathbf{C}'_{\text{right}}(\delta) \leq \mathbf{C}'_{\text{left}}(\delta) < \infty$  and the function  $\mathbf{C}(\delta)$  is increasing.*
3.  *$\mathbf{C}(\delta)$  is continuous at  $\delta = \delta_0$ ; therefore, it is continuous everywhere in  $[0, 1]$ .*

*If a function  $\mathbf{C} : [0, 1] \rightarrow [0, 1]$  satisfies these properties, there exist a measurable space  $\mathbf{X}$ , a finite set  $\mathbf{Y}$ , and a probability distribution  $P$  in  $\mathbf{X} \times \mathbf{Y}$  for which  $\mathbf{C}$  is the complementary success curve.*

**Proof sketch** The statement of the lemma follows from the fact that the complementary success curve  $\mathbf{S}$  can be obtained from the predictability distribution function  $F$  using these steps (labeling the horizontal and vertical axes as  $x$  and  $y$  respectively):

1. Invert  $F$ :  $F_1 := F^{-1}$ .
2. Integrate  $F_1$ :  $F_2(x) := \int_0^x F_1(t)dt$ .
3. Increase  $F_2$ :  $F_3(x) := F_2(x) + \delta_0$ , where  $\delta_0 := \int_0^1 F(x)dx$ .
4. Invert  $F_3$ :  $F_4 := F_3^{-1}$ .

It can be shown that  $\mathbf{C} = F_4$ , if we define  $g^{-1}(y) := \sup\{x : g(x) \leq y\}$  for non-decreasing  $g$  (so that  $g^{-1}$  is continuous on the right).  $\blacksquare$

Complement the protocol of §2 in which Nature plays  $P^\infty$  and Predictor plays  $\Gamma_{1-\delta}$  with the following variables:

$$\begin{aligned} \overline{\text{err}}_n &:= (P \times \mathbf{U})\{(x, y, \tau) : y \notin \Gamma_{1-\delta}(x_1, \tau_1, y_1, \dots, x_{n-1}, \tau_{n-1}, y_{n-1}, x, \tau)\}, \\ \overline{\text{unc}}_n &:= (P_{\mathbf{X}} \times \mathbf{U})\{(x, \tau) : |\Gamma_{1-\delta}(x_1, \tau_1, y_1, \dots, x_{n-1}, \tau_{n-1}, y_{n-1}, x, \tau)| > 1\}, \\ \overline{\text{emp}}_n &:= (P_{\mathbf{X}} \times \mathbf{U})\{(x, \tau) : |\Gamma_{1-\delta}(x_1, \tau_1, y_1, \dots, x_{n-1}, \tau_{n-1}, y_{n-1}, x, \tau)| = 0\}, \end{aligned}$$

and

$$\overline{\text{Err}}_n := \sum_{i=1}^n \overline{\text{err}}_i, \quad \overline{\text{Unc}}_n := \sum_{i=1}^n \overline{\text{unc}}_i, \quad \overline{\text{Emp}}_n := \sum_{i=1}^n \overline{\text{emp}}_i.$$

By the martingale strong law of large numbers it suffices to prove the proposition for these “predictable” versions of  $\text{Err}_n$ ,  $\text{Unc}_n$ , and  $\text{Emp}_n$ : indeed, since  $\text{Err}_n - \overline{\text{Err}}_n$  and  $\text{Unc}_n - \overline{\text{Unc}}_n$  are martingales (with increments bounded by 1 in absolute value) with respect to the filtration  $\mathcal{F}_n$ ,  $n = 0, 1, \dots$ , where each  $\mathcal{F}_n$  is generated by  $(x_1, \tau_1, y_1), \dots, (x_n, \tau_n, y_n)$ , we have

$$\lim_{n \rightarrow \infty} \frac{\text{Err}_n - \overline{\text{Err}}_n}{n} = 0 \quad \text{a.s.}$$

and

$$\lim_{n \rightarrow \infty} \frac{\text{Unc}_n - \overline{\text{Unc}}_n}{n} = 0 \quad \text{a.s.}$$

(see, e.g., [8], Theorem VII.5.4).

Without loss of generality we can assume that Nature's move  $\Gamma_n$  at trial  $n$  is  $\{\hat{y}(x_n)\}$  (where  $x \mapsto \hat{y}(x) \in \arg \max_y P(y|x)$  is a fixed "choice function") or the empty set  $\emptyset$  or the whole label space  $\mathbf{Y}$ . Furthermore, we can assume that

$$\overline{\text{unc}}_n = \mathbf{S}(\overline{\text{err}}_n), \quad \overline{\text{emp}}_n = \mathbf{C}(\overline{\text{err}}_n),$$

at every trial, since the optimal way (which is essentially the only optimal way) to spend the allowance of  $\overline{\text{err}}_n$  is to be certain on objects  $x$  with the largest (uppermost in Figure 2) representations  $F(f(x))$  and, if part of the allowance is still left when certainty is attained for all objects (i.e., if  $\overline{\text{err}}_n > \delta_0$ ), be empty on objects  $x$  with the smallest (lowermost in Figure 2) representations  $F(f(x))$ . (A formal argument for the statement about  $\mathbf{S}$  is given in [9].)

From the argument of [9] (Proof of Theorems 1 and 1r) it is clear that to achieve (8) the region predictor must satisfy

$$\begin{aligned} \delta < \delta_0 &\implies \limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (\overline{\text{err}}_i - \delta_0)^+ = 0 \\ \delta \geq \delta_0 &\implies \limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (\delta_0 - \overline{\text{err}}_i)^+ = 0. \end{aligned}$$

Using the fact that the complementary success curve  $\mathbf{C}$  is concave, increasing, and (uniformly) continuous for  $\delta \geq \delta_0$  (see Lemma 1), we obtain: if  $\delta < \delta_0$ ,

$$\begin{aligned} \frac{\overline{\text{Emp}}_n}{n} &= \frac{1}{n} \sum_{i=1}^n \overline{\text{emp}}_i = \frac{1}{n} \sum_{i=1}^n \mathbf{C}(\overline{\text{err}}_i) \\ &\leq \frac{1}{n} \mathbf{C}'_{\text{right}}(\delta_0) \sum_{i=1}^n (\overline{\text{err}}_i - \delta_0)^+ \rightarrow 0 \quad (n \rightarrow \infty); \end{aligned}$$

if  $\delta \geq \delta_0$ ,

$$\begin{aligned} \frac{\overline{\text{Emp}}_n}{n} &= \frac{1}{n} \sum_{i=1}^n \mathbf{C}(\overline{\text{err}}_i) = \frac{1}{n} \sum_{i=1}^n \mathbf{C}(\overline{\text{err}}_i \vee \delta_0) \\ &\leq \mathbf{C} \left( \frac{1}{n} \sum_{i=1}^n (\overline{\text{err}}_i \vee \delta_0) \right) = \mathbf{C} \left( \frac{1}{n} \sum_{i=1}^n \overline{\text{err}}_i + \frac{1}{n} \sum_{i=1}^n (\delta_0 - \overline{\text{err}}_i)^+ \right) \\ &\leq \mathbf{C} \left( \frac{1}{n} \sum_{i=1}^n \overline{\text{err}}_i \right) + o(1) \leq \mathbf{C}(\delta) + \epsilon, \end{aligned}$$

the last inequality holding almost surely for an arbitrary  $\epsilon > 0$  from some  $n$  on and  $\delta$  being the significance level used.

### A.3 Proof sketch of Proposition 5

Let us first modify and extend the notation  $P_n^\neq(y | x_i, \sigma_i)$  introduced in (6). Consider the sequence of extended examples  $w_i = (x_i, \tau'_i, y_i)$ ,  $i = 1, \dots, n$  ( $(x_i, y_i)$  are the first  $n$  examples chosen by Nature and  $\tau_i$  are the random numbers used by Predictor). We define the Nearest Neighbors approximations  $P_n(y | x, \sigma)$  to the conditional probabilities  $P(y | x)$  as follows: for every  $(x, \sigma, y) \in \tilde{\mathbf{Z}}$ ,

$$P_n(y | x, \sigma) := N(x, \sigma, y) / K_n,$$

where  $N(x, \sigma, y)$  is the number of  $i = 1, \dots, n$  such that  $(x_i, \tau'_i)$  is among the  $K_n$  nearest neighbors of  $(x, \sigma)$  and  $y_i = y$  (this time  $(x_i, \tau'_i)$  is not prevented from being counted as one of the  $K_n$  nearest neighbors of  $(x, \sigma)$  if  $(x_i, \tau'_i) = (x, \sigma)$ ). We define the empirical predictability function  $f_n$  by

$$f_n(x, \sigma) := \max_{y \in \mathbf{Y}} P_n(y | x, \sigma).$$

The proof will be based on the following version of a well-known fundamental result.

**Lemma 2** *Suppose  $K_n \rightarrow \infty$ ,  $K_n = o(n)$ , and  $\mathbf{Y} = \{0, 1\}$ . For any  $\epsilon > 0$  and large enough  $n$ ,*

$$\mathbb{P} \left\{ \int |P(1 | x) - P_n(1 | x, \sigma)| P_{\mathbf{X}}(dx) \mathbf{U}(d\sigma) > \epsilon \right\} \leq e^{-n\epsilon^2/40},$$

where the outermost probability distribution  $\mathbb{P}$  (essentially  $(P \times \mathbf{U})^\infty$ ) generates the extended examples  $(x_i, \tau_i, y_i)$ , which determine the empirical distributions  $P_n$ .

**Proof** This is almost a special case of Devroye et al.'s [3] Theorem 1. There is, however, an important difference between the way we break distance ties and the way Devroye et al. [3] do this: in [3], instead of our (2),

$$(|x_1 - x_3|, |\sigma_1 - \sigma_3|) < (|x_2 - x_3|, |\sigma_2 - \sigma_3|)$$

is used. (Our way of breaking ties better agrees with the lexicographic order on  $[0, 1]^2$ , which is useful in the proof of Proposition 2 and, less importantly, in the proof of Lemma 4.) It is easy to check that the proof given in [3] also works (and becomes simpler) for our way of breaking distance ties. ■

**Lemma 3** Suppose  $K_n \rightarrow \infty$  and  $K_n = o(n)$ . For any  $\epsilon > 0$  there exists an  $\epsilon^* > 0$  such that, for large enough  $n$ ,

$$\mathbb{P} \left\{ (P_{\mathbf{X}} \times \mathbf{U}) \left\{ (x, \sigma) : \max_{y \in \mathbf{Y}} |P_n(y|x, \sigma) - P(y|x)| > \epsilon \right\} > \epsilon \right\} \leq e^{-\epsilon^* n};$$

in particular,

$$\mathbb{P} \{ (P_{\mathbf{X}} \times \mathbf{U}) \{ (x, \sigma) : |f_n(x, \sigma) - f(x)| > \epsilon \} > \epsilon \} \leq e^{-\epsilon^* n}.$$

**Proof** We apply Lemma 2 to the binary classification problem obtained from our classification problem by replacing label  $y \in \mathbf{Y}$  with 1 and replacing all other labels with 0:

$$\mathbb{P} \left\{ \int |P(y|x) - P_n(y|x, \sigma)| P_{\mathbf{X}}(dx) \mathbf{U}(d\sigma) > \epsilon \right\} \leq e^{-n\epsilon^2/40}.$$

By Markov's inequality this implies

$$\mathbb{P} \{ (P_{\mathbf{X}} \times \mathbf{U}) \{ |P(y|x) - P_n(y|x, \sigma)| > \sqrt{\epsilon} \} > \sqrt{\epsilon} \} \leq e^{-n\epsilon^2/40},$$

which, in turn, implies

$$\mathbb{P} \left\{ (P_{\mathbf{X}} \times \mathbf{U}) \left\{ \max_{y \in \mathbf{Y}} |P(y|x) - P_n(y|x, \sigma)| > \sqrt{\epsilon} \right\} > |\mathbf{Y}| \sqrt{\epsilon} \right\} \leq e^{-n\epsilon^2/40}.$$

This completes the proof, since we can take the  $\epsilon$  in the last equation arbitrarily small as compared to the  $\epsilon$  in the statement of the lemma.  $\blacksquare$

We will use the shorthand “ $\forall^\infty n$ ” for “from some  $n$  on”.

**Lemma 4** Suppose  $K_n \rightarrow \infty$  and  $K_n = o(n)$ . For any  $\epsilon > 0$  there exists an  $\epsilon^* > 0$  such that, for large enough  $n$ ,

$$\mathbb{P} \left\{ \frac{\# \{ i : \max_y |P(y|x_i) - P_n^\#(y|x_i, \tau'_i)| > \epsilon \}}{n} > \epsilon \right\} \leq e^{-\epsilon^* n}.$$

In particular,

$$\forall^\infty n : \mathbb{P} \left\{ \frac{\# \{ i : |f(x_i) - f_n^\#(x_i, \tau'_i)| > \epsilon \}}{n} > \epsilon \right\} \leq e^{-\epsilon^* n}.$$

**Proof** Since

$$|P_n^{\neq}(y | x_i, \tau'_i) - P_n(y | x_i, \tau'_i)| \leq \frac{1}{K_n} = o(1),$$

we can, and will, ignore the upper indices  $\neq$  in the statement of the lemma.

Define

$$I_n(x, \sigma) := \begin{cases} 0 & \text{if } \max_y |P(y | x) - P_n(y | x, \sigma)| \leq \epsilon \\ 1 & \text{if } \max_y |P(y | x) - P_n(y | x, \sigma)| \geq 2\epsilon \\ (\max_y |P(y | x) - P_n(y | x, \sigma)| - \epsilon)/\epsilon & \text{otherwise} \end{cases}$$

(intuitively,  $I_n(x)$  is a “soft version” of  $\mathbb{I}_{\{\max_y |P(y|x) - P_n(y|x,\sigma)| > \epsilon\}}$ ).

The main tool in this proof (and several other proofs in this appendix) will be McDiarmid’s theorem (see, e.g., [4], Theorem 9.2). First we check the possibility of its application. If we replace an extended object  $(x_j, \tau'_j)$  by another extended object  $(x_j^*, \tau_j^*)$ , the expression

$$\sum_{i=1}^n I_n(x_i, \tau'_i)$$

will change as follows:

- the addend  $I_n(x_i, \tau'_i)$  for  $i = j$  changes by 1 at most;
- the addends  $I_n(x_i, \tau'_i)$  for  $i \neq j$  such that neither  $(x_j, \tau'_j)$  nor  $(x_j^*, \tau_j^*)$  are among the  $K_n$  nearest neighbors of  $(x_i, \tau'_i)$  do not change at all;
- the sum over the at most  $4K_n$  (see below) addends  $I_n(x_i, \tau'_i)$  for  $i \neq j$  such that either  $(x_j, \tau'_j)$  or  $(x_j^*, \tau_j^*)$  (or both) are among the  $K_n$  nearest neighbors of  $(x_i, \tau'_i)$  can change by at most

$$4K_n \frac{1}{\epsilon} \frac{1}{K_n} = \frac{4}{\epsilon}. \tag{11}$$

The left-hand side of (11) reflects the following facts: the change in  $P_n(y | x_i, \tau'_i)$  for  $i \neq j$  is at most  $1/K_n$ ; the number of  $i \neq j$  such that  $(x_j, \tau'_j)$  is among the  $K_n$  nearest neighbors of  $(x_i, \tau'_i)$  does not exceed  $2K_n$  (since the extended objects are linearly ordered and (2) is used for breaking distance ties); analogously, the number of  $i \neq j$  such that  $(x_j^*, \tau_j^*)$  is among the  $K_n$  nearest neighbors of  $(x_i, \tau'_i)$  does not exceed  $2K_n$ .

Therefore, by McDiarmid's theorem,

$$\begin{aligned} & \mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n I_n(x_i, \tau'_i) - \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n I_n(x_i, \tau'_i) \right) > \epsilon \right\} \\ & \leq \exp(-2\epsilon^2 n / (1 + 4/\epsilon)) = \exp\left(-\frac{2\epsilon^3}{4 + \epsilon} n\right). \end{aligned} \quad (12)$$

Next we find:

$$\begin{aligned} & \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n I_n(x_i, \tau'_i) \right) = \mathbb{E}(I_n(x_n, \tau'_n)) \leq \mathbb{E}(I_{n-1}(x_n, \tau'_n)) + o(1) \\ & \leq \mathbb{E}(P_{\mathbf{X}} \times \mathbf{U}) \{ (x, \sigma) : \max_y |P(y|x) - P_{n-1}(y|x, \sigma)| > \epsilon \} + o(1) \\ & \leq e^{-\epsilon^* n} + \epsilon + o(1) \leq 2\epsilon \end{aligned}$$

(the penultimate inequality follows from Lemma 3) from some  $n$  on. In combination with (12) this implies

$$\forall^\infty n : \mathbb{P} \left\{ \frac{1}{n} \sum_{i=1}^n I_n(x_i, \tau'_i) > 3\epsilon \right\} \leq \exp\left(-\frac{2\epsilon^3}{4 + \epsilon} n\right),$$

in particular

$$\mathbb{P} \left\{ \frac{\#\{i : \max_y |P(y|x_i) - P_n(y|x_i, \tau'_i)| \geq 2\epsilon\}}{n} > 3\epsilon \right\} \leq \exp\left(-\frac{2\epsilon^3}{4 + \epsilon} n\right).$$

Replacing  $3\epsilon$  by  $\epsilon$ , we obtain that, from some  $n$  on,

$$\mathbb{P} \left\{ \frac{\#\{i : \max_y |P(y|x_i) - P_n(y|x_i, \tau'_i)| > \epsilon\}}{n} > \epsilon \right\} \leq \exp\left(-\frac{2(\epsilon/3)^3}{4 + \epsilon/3} n\right),$$

which completes the proof. ■

We say that an extended example  $(x_i, \tau_i, y_i)$ ,  $i = 1, \dots, n$ , is *n-strange* if  $y_i \neq \hat{y}_n(x_i, \tau'_i)$ ; otherwise,  $(x_i, \tau_i, y_i)$  will be called *n-ordinary*. For concreteness, let us break ties in the expression “the  $k$  largest” by regarding elements corresponding to larger indices as larger for this purpose (although with probability one, ties will never happen in our case).

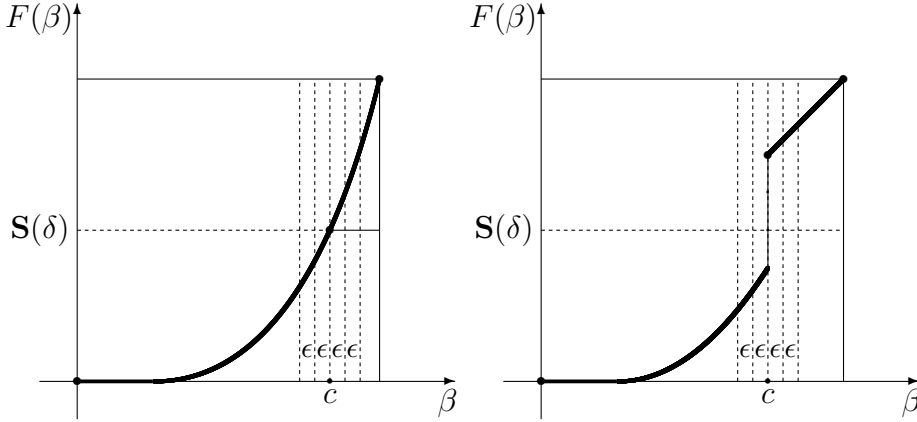


Figure 3: Cases  $F(c) = \mathbf{S}(\delta)$  (left) and  $F(c) > \mathbf{S}(\delta)$  (right)

**Lemma 5** *Suppose (5) is satisfied and  $\delta \leq \delta_0$ . With probability one, the  $\lfloor (1 - \mathbf{S}(\delta))n \rfloor$  extended examples with the largest (in the sense of the lexicographic order)  $(f_n^\#(x_i, \tau_i'), \tau_i'')$  among  $(x_1, \tau_1, y_1), \dots, (x_n, \tau_n, y_n)$  contain at most  $n\delta + o(n)$   $n$ -strange extended examples as  $n \rightarrow \infty$ .*

**Proof sketch** Define

$$c := \sup\{\beta : F(\beta) \leq \mathbf{S}(\delta)\}.$$

It is clear that  $0 < c < 1$ . Our proof will work both in the case where  $F(c) = \mathbf{S}(\delta)$  and in the case where  $F(c) > \mathbf{S}(\delta)$ , as illustrated in Figure 3.

Let  $\epsilon > 0$  be a small constant (we will let  $\epsilon \rightarrow 0$  eventually). Define a “threshold”  $(c'_n, c''_n) \in [0, 1]^2$  requiring that

$$\mathbb{P}\{f(x_n) = c, (f_{n-1}(x_n, \tau'_n), \tau''_n) > (c'_n, c''_n)\} = F(c) - \mathbf{S}(\delta) - \epsilon \quad (13)$$

if  $F(c) > \mathbf{S}(\delta)$ ; we assume that  $\epsilon$  is small enough for

$$2\epsilon < F(c) - \mathbf{S}(\delta) \quad (14)$$

to hold (among other things this will ensure the validity of the definition (13)). If  $F(c) = \mathbf{S}(\delta)$ , we set  $(c', c'') := (c + \epsilon, 0)$ ; in any case, we will have

$$\mathbb{P}\{f(x_n) = c, (f_{n-1}(x_n, \tau'_n), \tau''_n) > (c'_n, c''_n)\} \geq F(c) - \mathbf{S}(\delta) - \epsilon. \quad (15)$$



Let us say that an extended example  $(x_i, \tau_i, y_i)$  is *above the threshold* if

$$(f_n^\neq(x_i, \tau_i'), \tau_i'') > (c'_n, c''_n);$$

otherwise, we say it is *below the threshold*. Divide the first  $n$  extended examples  $(x_i, \tau_i, y_i)$ ,  $i = 1, \dots, n$ , into five classes:

**Class I:** Those satisfying  $f(x_i) \leq c - 2\epsilon$ .

**Class II:** Those that satisfy  $f(x_i) = c$  and are below the threshold.

**Class III:** Those satisfying  $c - 2\epsilon < f(x_i) \leq c + 2\epsilon$  but not  $f(x_i) = c$ .

**Class IV:** Those that satisfy  $f(x_i) = c$  and are above the threshold.

**Class V:** Those satisfying  $f(x_i) > c + 2\epsilon$ .

First we explain the general idea of the proof. The threshold  $(c', c'')$  was chosen so that approximately  $\lfloor (1 - \mathbf{S}(\delta))n \rfloor$  of the available extended examples will be above the threshold. Because of this, the extended examples above the threshold will essentially be the  $\lfloor (1 - \mathbf{S}(\delta))n \rfloor$  extended examples with the largest  $(f_n^\neq(x_i, \tau_i'), \tau_i'')$  referred to in Lemma 5. For each of our five classes we will be interested in the following questions:

- How many extended examples are there in the class?
- How many of those are above the threshold?
- How many of those above the threshold are  $n$ -strange?

If the sum of the answers to the last question does not exceed  $n\delta$  by too much, we are done.

With this plan in mind, we start the formal proof. (Of course, we will not be following the plan literally: for example, if a class is very small, we do not need to answer the second and third questions.) The first step is to show that

$$c - \epsilon \leq c'_n \leq c + \epsilon \tag{16}$$

from some  $n$  on; this will ensure that the classes are conveniently separated from each other. We only need to consider the case  $F(c) > \mathbf{S}(\delta)$ . The inequality  $c'_n \leq c + \epsilon$  follows from

$$\forall^\infty n : \mathbb{P} \{f(x_n) = c, f_{n-1}(x_n, \tau'_n) > c + \epsilon\} < \epsilon < F(c) - \mathbf{S}(\delta) - \epsilon$$

(combine Lemma 3 with (14)). The inequality  $c - \epsilon \leq c'_n$  follows from

$$\begin{aligned} \forall^\infty n : \quad & \mathbb{P}\{f(x_n) = c, f_{n-1}(x_n, \tau'_n) \geq c - \epsilon\} \\ &= \mathbb{P}\{f(x_n) = c\} - \mathbb{P}\{f(x_n) = c, f_{n-1}(x_n, \tau'_n) < c - \epsilon\} \\ &> F(c) - F(c-) - \epsilon \geq F(c) - \mathbf{S}(\delta) - \epsilon. \end{aligned}$$

Now we are ready to analyze the composition of our five classes. Among the Class I extended examples at most

$$\epsilon n \tag{17}$$

will be above the threshold from some  $n$  on almost surely (by Lemma 4 and the Borel-Cantelli lemma). None of the Class II extended examples will be above the threshold, by definition. The fraction of Class III extended examples among the first  $n$  extended examples will tend to

$$F(c + 2\epsilon) - F(c) + F(c-) - F(c - 2\epsilon) \tag{18}$$

as  $n \rightarrow \infty$  almost surely.

To estimate the number  $N_n^{\text{IV}}$  of Class IV extended examples among the first  $n$  extended examples, we use McDiarmid's theorem. If one extended example is replaced by another,  $N_n^{\text{IV}}$  will change by at most  $2K_n + 1$  (since this extended example can affect  $f_n^\neq(x_i, \tau'_i)$  for at most  $2K_n$  other extended examples  $(x_i, \tau_i, y_i)$ ). Therefore,

$$\mathbb{P}\left\{\left|\frac{1}{n}N_n^{\text{IV}} - \frac{1}{n}\mathbb{E}N_n^{\text{IV}}\right| \geq \epsilon\right\} \leq 2e^{-2\epsilon^2n/(2K_n+1)},$$

the assumption  $K_n = o(n/\ln n)$  and the Borel-Cantelli lemma imply that

$$\left|\frac{1}{n}N_n^{\text{IV}} - \frac{1}{n}\mathbb{E}N_n^{\text{IV}}\right| < \epsilon$$

from some  $n$  on almost surely. Since

$$\frac{1}{n}\mathbb{E}N_n^{\text{IV}} = \mathbb{P}\{f(x_n) = c, (f_{n-1}(x_n, \tau'_n), \tau''_n) > (c'_n, c''_n)\} \geq F(c) - \mathbf{S}(\delta) - \epsilon$$

(see (15)), we have

$$N_n^{\text{IV}} > (F(c) - \mathbf{S}(\delta) - 2\epsilon)n \tag{19}$$

from some  $n$  on almost surely. Of course, all these examples are above the threshold.

Now we estimate the number  $N_n^{\text{IV, str}}$  of  $n$ -strange extended examples of Class IV. Again McDiarmid's theorem implies that

$$\left| \frac{1}{n} N_n^{\text{IV, str}} - \frac{1}{n} \mathbb{E} N_n^{\text{IV, str}} \right| < \epsilon$$

from some  $n$  on almost surely. Now, from some  $n$  on,

$$\begin{aligned} \frac{1}{n} \mathbb{E} N_n^{\text{IV, str}} &= \mathbb{P} \{ f(x_n) = c, (f_{n-1}(x_n, \tau'_n), \tau''_n) > (c'_n, c''_n), \hat{y}_n(x_n, \tau'_n) \neq y_n \} \\ &= \mathbb{E} \left( (1 - P(\hat{y}_n(x_n, \tau'_n) | x_n)) \mathbb{I}_{\{f(x_n)=c, (f_{n-1}(x_n, \tau'_n), \tau''_n) > (c'_n, c''_n)\}} \right) \\ &\leq e^{-\epsilon^* n} + \epsilon + (1 - c + \epsilon) \\ &\quad \times \mathbb{P} \{ f(x_n) = c, (f_{n-1}(x_n, \tau'_n), \tau''_n) > (c'_n, c''_n) \} \\ &= e^{-\epsilon^* n} + \epsilon + (1 - c + \epsilon)(F(c) - \mathbf{S}(\delta) - \epsilon) \tag{20} \\ &\leq (F(c) - \mathbf{S}(\delta))(1 - c) + 3\epsilon \tag{21} \end{aligned}$$

(the first inequality follows from Lemma 3) in the case  $F(c) > \mathbf{S}(\delta)$ . If  $F(c) = \mathbf{S}(\delta)$ , the lines (20) and (21) of this chain have to be changed to

$$\begin{aligned} &\leq e^{-\epsilon^* n} + \epsilon + (1 - c + \epsilon) (e^{-\epsilon^* n} + \epsilon) \\ &< 3\epsilon, \end{aligned}$$

but the inequality between the extreme terms of the chain still holds. Therefore, the number of  $n$ -strange Class IV extended examples does not exceed

$$((F(c) - \mathbf{S}(\delta))(1 - c) + 4\epsilon) n \tag{22}$$

from some  $n$  on almost surely.

By the Borel strong law of large numbers, the fraction of Class V extended examples among the first  $n$  extended examples will tend to

$$1 - F(c + 2\epsilon) \tag{23}$$

as  $n \rightarrow \infty$  almost surely. By Lemma 4, the Borel-Cantelli lemma, and (16), almost surely from some  $n$  on at least

$$(1 - F(c + 2\epsilon) - 2\epsilon)n \tag{24}$$

extended examples in Class V will be above the threshold.

Finally, we estimate the number  $N_n^{\text{V, str}}$  of  $n$ -strange extended examples of Class V among the first  $n$  extended examples. By McDiarmid's theorem,

$$\left| \frac{1}{n} N_n^{\text{V, str}} - \frac{1}{n} \mathbb{E} N_n^{\text{V, str}} \right| < \epsilon$$

from some  $n$  on almost surely. Now

$$\begin{aligned} \frac{1}{n} \mathbb{E} N_n^{\text{V, str}} &= \mathbb{P} \{f(x_n) > c + 2\epsilon, \hat{y}_n(x_n, \tau'_n) \neq y_n\} \\ &= \mathbb{E} \left( (1 - P(\hat{y}_n(x_n, \tau'_n) | x_n)) \mathbb{I}_{\{f(x_n) > c + 2\epsilon\}} \right) \\ &\leq e^{-\epsilon^* n} + \epsilon + \mathbb{E} \left( (1 - f(x_n) + \epsilon) \mathbb{I}_{\{f(x_n) > c + 2\epsilon\}} \right) \\ &\leq e^{-\epsilon^* n} + 2\epsilon + \mathbb{E} \left( (1 - f(x_n)) \mathbb{I}_{\{f(x_n) > c + 2\epsilon\}} \right) \\ &= e^{-\epsilon^* n} + 2\epsilon + \int_0^1 (F(\beta) - F(c + 2\epsilon))^+ d\beta \\ &< \int_0^1 (F(\beta) - F(c))^+ d\beta + 3\epsilon \end{aligned}$$

from some  $n$  on (the first inequality follows from Lemma 3). Therefore,

$$\frac{1}{n} N_n^{\text{V, str}} < \int_0^1 (F(\beta) - F(c))^+ d\beta + 4\epsilon \quad (25)$$

from some  $n$  on almost surely.

Summarizing, we can see that the total number of extended examples above the threshold among the first  $n$  extended examples will be at least

$$\begin{aligned} &(F(c) - \mathbf{S}(\delta) - 2\epsilon + 1 - F(c + 2\epsilon) - 2\epsilon) n \\ &= (1 - \mathbf{S}(\delta) + F(c) - F(c + 2\epsilon) - 4\epsilon) n \end{aligned} \quad (26)$$

(see (19) and (24)) from some  $n$  on almost surely. The number of  $n$ -strange extended examples among them will not exceed

$$\begin{aligned} &\left( \epsilon + F(c + 2\epsilon) - F(c) + F(c-) - F(c - 2\epsilon) + \epsilon \right. \\ &+ (F(c) - \mathbf{S}(\delta))(1 - c) + 4\epsilon + \int_0^1 (F(\beta) - F(c))^+ d\beta + 4\epsilon \left. \right) n \\ &= \left( F(c + 2\epsilon) - F(c) + F(c-) - F(c - 2\epsilon) \right. \\ &+ (F(c) - \mathbf{S}(\delta))(1 - c) + \int_0^1 (F(\beta) - F(c))^+ d\beta + 10\epsilon \left. \right) n \end{aligned} \quad (27)$$

(see (17), (18), (22), and (25)) from some  $n$  on almost surely. Combining (26) and (27), we can see that the number of  $n$ -strange extended examples among the  $\lfloor (1 - \mathbf{S}(\delta))n \rfloor$  extended examples with the largest  $(f_n^\neq(x_i, \tau_i'), \tau_i'')$  does not exceed

$$\begin{aligned} & \left( F(c + 2\epsilon) - F(c) + F(c-) - F(c - 2\epsilon) + (F(c) - \mathbf{S}(\delta))(1 - c) \right. \\ & \quad \left. + \int_0^1 (F(\beta) - F(c))^+ d\beta + 10\epsilon \right) n + (F(c + 2\epsilon) - F(c) + 4\epsilon) n \\ = & \left( 2(F(c + 2\epsilon) - F(c)) + (F(c-) - F(c - 2\epsilon)) + (F(c) - \mathbf{S}(\delta))(1 - c) \right. \\ & \quad \left. + \int_0^1 (F(\beta) - F(c))^+ d\beta + 14\epsilon \right) n \end{aligned}$$

from some  $n$  on almost surely. Since  $\epsilon$  can be arbitrarily small, the coefficient in front of  $n$  in the last expression can be made arbitrarily close to

$$(F(c) - \mathbf{S}(\delta))(1 - c) + \int_0^1 (F(\beta) - F(c))^+ d\beta = \int_0^1 (F(\beta) - \mathbf{S}(\delta))^+ d\beta = \delta,$$

which completes the proof.  $\blacksquare$

**Lemma 6** *Suppose (5) is satisfied. The fraction of  $n$ -strange extended examples among the first  $n$  extended examples  $(x_i, \tau_i, y_i)$  approaches  $\delta_0$  asymptotically with probability one.*

**Proof sketch** The lemma is not difficult to prove using McDiarmid's theorem and the fact that, by Lemma 4,  $P(\hat{y}_n(x_i, \tau_i') | x_i)$  will typically differ little from  $f(x_i)$ . Notice, however, that the part that we really need in this paper (that the fraction of  $n$ -strange extended examples does not exceed  $\delta_0 + o(1)$  as  $n \rightarrow \infty$  with probability one) is just a special case of Lemma 5, corresponding to  $\delta = \delta_0$ .  $\blacksquare$

**Lemma 7** *Suppose (5) is satisfied and  $\delta > \delta_0$ . The fraction of  $n$ -ordinary extended examples among the  $\lfloor \mathbf{C}(\delta)n \rfloor$  extended examples  $(x_i, \tau_i, y_i)$ ,  $i = 1, \dots, n$ , with the lowest  $(f_n^\neq(x_i, \tau_i'), \tau_i'')$  does not exceed  $\delta - \delta_0 + o(1)$  as  $n \rightarrow \infty$  with probability one.*

Lemma 7 can be proven analogously to Lemma 5.

**Lemma 8** Let  $\mathcal{F}_1 \supseteq \mathcal{F}_2 \supseteq \dots$  be a decreasing sequence of  $\sigma$ -algebras and  $\xi_1, \xi_2, \dots$  be a bounded adapted (in the sense that  $\xi_n$  is  $\mathcal{F}_n$ -measurable for all  $n$ ) sequence of random variables such that

$$\limsup_{n \rightarrow \infty} \mathbb{E}(\xi_n | \mathcal{F}_{n+1}) \leq 0 \quad a.s.$$

Then

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \xi_i \leq 0 \quad a.s.$$

**Proof** Replacing, if necessary,  $\xi_n$  by  $\xi_n - \mathbb{E}(\xi_n | \mathcal{F}_{n+1})$ , we reduce our task to the following special case (a reverse Borel strong law of large numbers): if  $\xi_1, \xi_2, \dots$  is a bounded *reverse martingale difference*, in the sense of being adapted and satisfying  $\forall n : \mathbb{E}(\xi_n | \mathcal{F}_{n+1}) = 0$ , then

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \xi_i = 0 \quad a.s. \quad (28)$$

Fix a bounded reverse martingale difference  $\xi_1, \xi_2, \dots$ ; our goal is to prove (28). By the martingale version of Hoeffding's inequality ([4], Theorem 9.1) applied to the martingale difference  $(\xi_i, \mathcal{F}_i)$ ,  $i = n, \dots, 1$ ,

$$\mathbb{P} \left\{ \left| \frac{1}{n} \sum_{i=1}^n \xi_i \right| \geq \epsilon \right\} \leq 2e^{-2\epsilon^2 n / (2C)}, \quad (29)$$

where  $C$  is an upper bound on  $\sup_n |\xi_n|$ . Combined with the Borel-Cantelli-Lévy lemma, (29) implies (28).  $\blacksquare$

Now we can sketch the proof of Proposition 5. Define  $\mathcal{F}_n$ ,  $n = 1, 2, \dots$ , to be the  $\sigma$ -algebra on  $\tilde{\mathbf{Z}}^\infty$  generated by the multiset of the first  $n - 1$  extended examples  $(x_i, \tau_i, y_i)$ ,  $i = 1, \dots, n - 1$ , and the sequence of extended examples  $(x_i, \tau_i, y_i)$ ,  $i = n, n + 1, \dots$  (starting from the  $n$ th extended example).

Suppose first that  $\delta < \delta_0$ . Consider the  $\lfloor (1 - \mathbf{S}(\delta - \epsilon))n \rfloor$  extended examples with the largest  $(f_n^\#(x_i, \tau_i'), \tau_i'')$  among  $(x_1, \tau_1, y_1), \dots, (x_n, \tau_n, y_n)$ , where  $\epsilon \in (0, \delta)$  is a small constant. Let us show that each of these examples will be predicted with certainty from the other extended examples in the sequence  $(x_1, \tau_1, y_1), \dots, (x_n, \tau_n, y_n)$ , from some  $n$  on. We will be assuming  $n$  large enough.

Let  $(x_k, \tau_k, y_k)$  be the extended example with the  $(\lfloor(\delta - \epsilon/2)n\rfloor + 1)$ th largest (in the sense of the lexicographic order)  $(f_n^\neq(x_i, \tau_i'), \tau_i'')$  among all  $n$ -strange extended examples  $(x_i, \tau_i, y_i)$ ,  $i = 1, \dots, n$ . (The simple case where  $(x_k, \tau_k, y_k)$  does not exist needs to be considered separately.) Let  $(x_j, \tau_j, y_j)$  be one of the  $\lfloor(1 - \mathbf{S}(\delta - \epsilon))n\rfloor$  extended examples with the largest  $(f_n^\neq(x_i, \tau_i'), \tau_i'')$  and let  $y \in \mathbf{Y}$  be a label different from  $\hat{y}_n(x_j, \tau_j')$ . It suffices to prove that

$$\tau_j'' \geq \frac{\#\{i = 1, \dots, n : \alpha_i^y \geq \alpha_j^y\} - n\delta}{\#\{i = 1, \dots, n : \alpha_i^y = \alpha_j^y\}} \quad (30)$$

(cf. (4) on p. 7), where all  $\alpha^y$  are computed as  $\alpha$  in (7) from the sequence  $(x_1, \tau_1, y_1), \dots, (x_n, \tau_n, y_n)$  with  $y_j$  replaced by  $y$ . It will be more convenient to write (30) in the form

$$\#\{i : \alpha_i^y > \alpha_j^y\} + (1 - \tau_j'')\#\{i : \alpha_i^y = \alpha_j^y\} \leq n\delta.$$

Since  $\alpha_j^y = f_n^\neq(x_j, \tau_j')$  and  $\alpha_i^y \neq \alpha_i$  for at most  $2K_n + 1$  values of  $i$  (indeed, changing  $y_j$  will affect at most  $2K_n + 1$   $\alpha$ s), it suffices to prove

$$\#\{i : \alpha_i > f_n^\neq(x_j, \tau_j')\} + (1 - \tau_j'')\#\{i : \alpha_i = f_n^\neq(x_j, \tau_j')\} \leq n(\delta - \epsilon^*), \quad (31)$$

where  $\epsilon^* \ll \epsilon$  is a positive constant.

Since  $(f_n^\neq(x_j, \tau_j'), \tau_j'') \geq (\alpha_k, \tau_k'')$  (indeed, by Lemma 5, there are less than  $(\delta - \epsilon/2)n$   $n$ -strange extended examples among the  $\lfloor(1 - \mathbf{S}(\delta - \epsilon))n\rfloor$  extended examples with the largest  $(f_n^\neq(x_i, \tau_i'), \tau_i'')$ ), (31) will follow from

$$\#\{i : \alpha_i > \alpha_k\} + (1 - \tau_k'')\#\{i : \alpha_i = \alpha_k\} \leq n(\delta - \epsilon^*). \quad (32)$$

If  $\#\{i : \alpha_i = \alpha_k\} \leq \frac{\epsilon}{3}n$ , the left-hand side of (32) does not exceed

$$\left(\delta - \frac{\epsilon}{2}\right)n + \frac{\epsilon}{3}n < n(\delta - \epsilon^*),$$

so we can, and will, assume without loss of generality that

$$\#\{i : \alpha_i = \alpha_k\} > \frac{\epsilon}{3}n. \quad (33)$$

Since  $\tau_i''$  for the extended examples satisfying  $\alpha_i = \alpha_k$  are output according to the uniform distribution  $\mathbf{U}$ , the expected value of  $1 - \tau_k''$  is about

$$\frac{(\delta - \epsilon/2)n - \#\{i : \alpha_i > \alpha_k\}}{\#\{i : \alpha_i = \alpha_k\}},$$

and so by Hoeffding's inequality and the Borel-Cantelli lemma we will have (from some  $n$  on)

$$1 - \tau_k'' \leq \frac{(\delta - \epsilon/2)n - \#\{i: \alpha_i > \alpha_k\}}{\#\{i: \alpha_i = \alpha_k\}} + \epsilon^* \quad (34)$$

(remember (33)). Equation (32) will hold because its left-hand side can be transformed using (34) as

$$\begin{aligned} \#\{i: \alpha_i > \alpha_k\} + (1 - \tau_k'')\#\{i: \alpha_i = \alpha_k\} &\leq (\delta - \epsilon/2)n + \epsilon^* \#\{i: \alpha_i = \alpha_k\} \\ &\leq (\delta - \epsilon/2 + \epsilon^*)n \leq (\delta - \epsilon^*)n. \end{aligned}$$

The assertion we have just proved means that, almost surely from some  $n$  on,

$$\mathbb{P}(\{\text{unc}_n = 0\} | \mathcal{F}_{n+1}) \geq \frac{\lfloor (1 - \mathbf{S}(\delta - \epsilon))n \rfloor}{n} \geq 1 - \mathbf{S}(\delta - \epsilon) - \frac{1}{n}.$$

Since  $\epsilon$  can be arbitrarily small and  $\mathbf{S}$  is continuous ([9], Lemma 2), this implies

$$\limsup_{n \rightarrow \infty} \mathbb{E}(\text{unc}_n | \mathcal{F}_{n+1}) \leq \mathbf{S}(\delta) \quad \text{a.s.}$$

By Lemma 8 this implies, in turn,

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \text{unc}_i \leq \mathbf{S}(\delta) \quad \text{a.s.},$$

which coincides with (9).

If  $\delta \geq \delta_0$ , Lemma 6 implies that

$$\lim_{n \rightarrow \infty} \mathbb{E}(\text{unc}_n | \mathcal{F}_{n+1}) = 0 \quad \text{a.s.};$$

in combination with Lemma 8 this again implies (9).

Inequality (10) is treated in a similar way. Lemmas 6 and 7 imply that

$$\liminf_{n \rightarrow \infty} \mathbb{E}(\text{emp}_n | \mathcal{F}_{n+1}) \geq \mathbf{C}(\delta) \quad \text{a.s.}$$

(this inequality is vacuously true when  $\delta \leq \delta_0$ ). Another application of Lemma 8 gives

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \text{emp}_i \geq \mathbf{C}(\delta) \quad \text{a.s.},$$

i.e., (10).



**Remark** The derivation of Proposition 5 from Lemmas 5–8 would be very simple if we defined the individual strangeness measure by, say,

$$\alpha_i := \begin{cases} (-f_n^\#(x_i, \sigma_i), \sigma_i) & \text{if } y_i = \hat{y}_n(x_i, \sigma_i) \\ (f_n^\#(x_i, \sigma_i), \sigma_i) & \text{otherwise} \end{cases}$$

(with the lexicographic order on  $\alpha$ 's) instead of (7) (in which case the denominator of (4) would be 1 almost surely). Our definition (7), however, is simpler and, most importantly, facilitates the proof of Proposition 2.

## On-line Compression Modelling Project Working Papers

1. *On-line confidence machines are well-calibrated*, by Vladimir Vovk, April 2002.
2. *Asymptotic optimality of Transductive Confidence Machine*, by Vladimir Vovk, May 2002.
3. *Universal well-calibrated algorithm for on-line classification*, by Vladimir Vovk, November 2002.
4. *Mondrian Confidence Machine*, by Vladimir Vovk, David Lindsay, Ilia Nouretdinov and Alex Gammerman, March 2003.
5. *Testing exchangeability on-line*, by Vladimir Vovk, Ilia Nouretdinov and Alex Gammerman, February 2003.
6. *Criterion of calibration for Transductive Confidence Machine with limited feedback*, by Ilia Nouretdinov and Vladimir Vovk, April 2003.
7. *Online region prediction with real teachers*, by Daniil Ryabko, Vladimir Vovk and Alex Gammerman, March 2003.
8. *Well-calibrated predictions from on-line compression models*, by Vladimir Vovk, April 2003.