

Asymptotic optimality of Transductive Confidence Machine

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 #2

May 27, 2002

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

Abstract

Transductive Confidence Machine (TCM) is a way of converting standard machine-learning algorithms into algorithms that output predictive regions rather than point predictions. It has been shown recently that TCM is well-calibrated when used in the on-line mode: at any confidence level $1 - \delta$, the long-run relative frequency of errors is guaranteed not to exceed δ provided the examples are generated independently from the same probability distribution P . Therefore, the number of “uncertain” predictive regions (i.e., those containing more than one label) becomes the sole measure of performance. The main result of this paper is that for any probability distribution P (assumed to generate the examples), it is possible to construct a TCM (guaranteed to be well-calibrated even if the assumption is wrong) that performs asymptotically as well as the best region predictor under P .

Contents

1	Region Predictors	1
2	Well-Calibrated and Asymptotically Optimal Region Predictors	3
3	Transductive Confidence Machine	7
4	Asymptotically Optimal TCM	8
5	Randomised Region Predictors	8
6	Conclusion	10

1 Region Predictors

The notion of TCM was introduced in [5] and [8] (our exposition is, however, self-contained). Before we define TCM (in §3) we discuss general properties of region predictors.

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 σ -algebra of all subsets). The protocol includes variables Err_n (the total number of errors made up to and including trial n) and err_n (the binary variable showing if an error is made at trial n); it also includes analogous variables Unc_n and unc_n for uncertain predictions:

```

Err0 := 0;  Unc0 := 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 ;
END FOR.

```

We will use the notation $\mathbf{Z} := \mathbf{X} \times \mathbf{Y}$ for the *example space*; Γ_n will be called a *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^∞). This is Nature's randomised strategy.

A *region predictor* is a family (indexed by $\gamma \in [0, 1]$) of Predictor's strategies Γ_γ such that $\Gamma_\gamma(x_1, y_1, \dots, x_{n-1}, y_{n-1}, x_n)$ is a measurable function of Nature's moves and $\Gamma_{\gamma_1}(x_1, y_1, \dots, x_{n-1}, y_{n-1}, x_n) \subseteq \Gamma_{\gamma_2}(x_1, y_1, \dots, x_{n-1}, y_{n-1}, x_n)$ when $\gamma_1 \leq \gamma_2$. Since we are interested in prediction with confidence, the 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 95% and 99%); the complementary value δ is called the *significance level*.

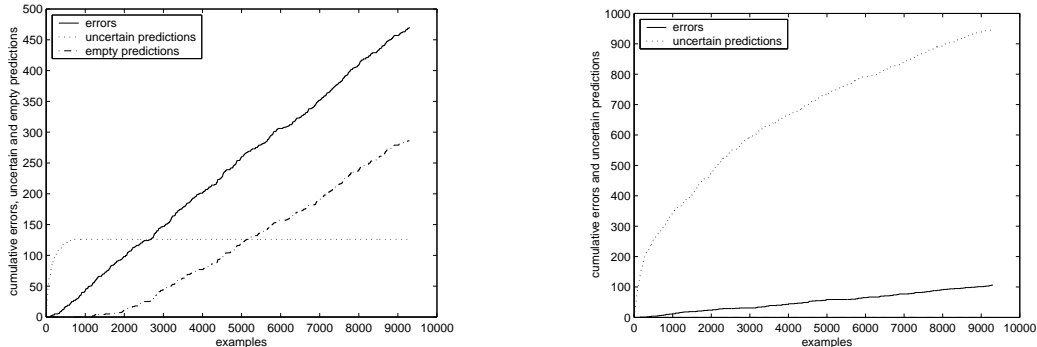


Figure 1: On-line performance of the Nearest Neighbour TCM on the USPS data set (9298 hand-written digits, randomly permuted) for the confidence level 95% (**left**) and 99% (**right**). The solid line shows the cumulative number of errors, dotted the cumulative number of uncertain predictions, and dashdot the cumulative number of empty predictions (inevitably leading to an error). For 99%, the dashdot line coincides with the horizontal axis (there are no empty predictions) and so is invisible. This and following figures are not significantly affected by statistical variation (due to the random choice of the permutation of the data set).

To provide the reader with an intuition about region prediction, we present results for Predictor’s particular strategy (“1-Nearest Neighbour TCM”; for details, see [10]) on the USPS data set (as described in [7], §12.2, but randomly permuted). Figure 1 shows the cumulative number of errors Err_n plotted against $n = 0, 1, \dots, 9298$ (solid line), the cumulative number Unc_n of uncertain predictions Γ_n , and that of empty predictions (for which $|\Gamma_n| = 0$). Figure 2 (left) gives the *empirical calibration curve*

$$\delta \mapsto \text{Err}_N(\text{USPS}, \Gamma_{1-\delta})/N$$

and the *empirical performance curve*

$$\delta \mapsto \text{Unc}_N(\text{USPS}, \Gamma_{1-\delta})/N$$

for this region predictor; we use the strategies followed by Nature (the randomly permuted USPS data set) and Predictor (corresponding to significance level δ) as arguments for Err and Unc ; $N = 9298$ is the size of the USPS data set.

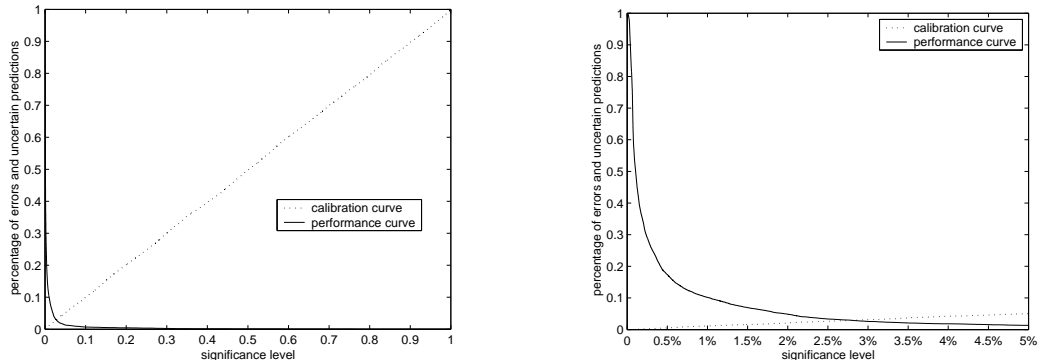


Figure 2: The empirical calibration and performance curves for the Nearest Neighbour TCM on the USPS data set (**left**); their left edges stretched horizontally (**right**).

2 Well-Calibrated and Asymptotically Optimal Region Predictors

Suppose we know the true distribution P in \mathbf{Z} generating the examples. In this section we will construct a region predictor optimal under P ; we will often omit P from our notation.

Let $P_{\mathbf{X}}$ be the marginal distribution of P in \mathbf{X} (i.e., $P_{\mathbf{X}}(E) := P(E \times \mathbf{Y})$) and $P_{\mathbf{Y}|\mathbf{X}}(y|x)$ be 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 subindices \mathbf{X} and $\mathbf{Y}|\mathbf{X}$.

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 3 (left); the graph of F is the thick line, and the unit box is also shown. The intuition behind some constructions in this paper will become clearer if the case of finite \mathbf{X} with equiprobable objects is considered first; see Figure 3, right.

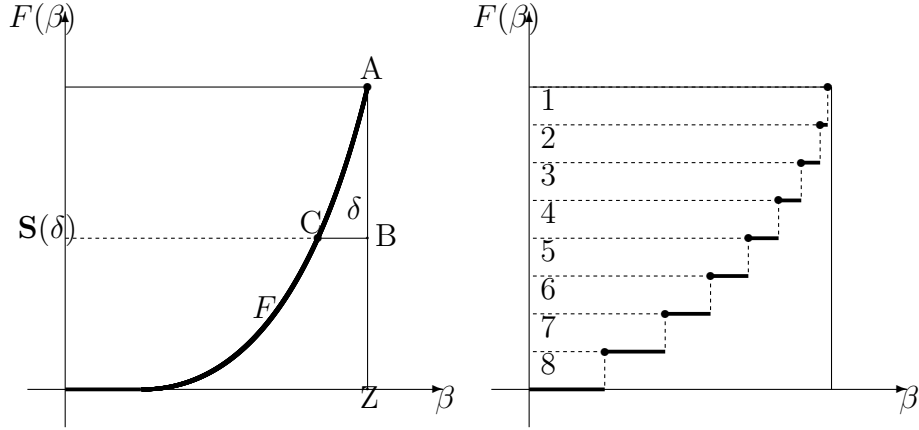


Figure 3: **Left:** The predictability distribution function F . The function F is non-decreasing, continuous on the right, and $F(1/|\mathbf{Y}|) = 0$. For a possibly more realistic example of a predictability distribution function, see Figure 6. **Right:** The predictability distribution function (thick line) in the case where the object space \mathbf{X} is finite and all objects $x \in \mathbf{X}$ have the same probability. The objects are numbered, from 1 to 8 in this case, in the order of decreasing predictability.

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

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

where t^+ stands for $\max(t, 0)$; the function \mathbf{S} is also of the type $[0, 1] \rightarrow [0, 1]$. (Why the terminology introduced here and below is natural will become clear from Theorems 1 and 2.) Geometrically, \mathbf{S} is defined from the graph of F as follows (see Figure 3, left): move the point B from A to Z until the area of the curvilinear triangle ABC becomes δ (assuming this area does become δ eventually, i.e., δ is not too large); the ordinate of B is then $\mathbf{S}(\delta)$. The intuition in the case of finite \mathbf{X} (Figure 3, right) is that $1 - \mathbf{S}(\delta)$ is the maximum fraction of objects that are “easily predictable” in the sense that their cumulative lack of predictability does not exceed δ (where the lack of predictability $1 - f(x)$ of each object is taken with the weight $1/|\mathbf{X}|$). Notice that the value $\mathbf{S}(\delta)$ in fact satisfies the equality

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

provided δ does not exceed the *critical significance level*

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

(the area under the thick curve in Figure 3, left; we will later see that this coincides with what is sometimes called *Bayes error* or *Bayes risk*—see, e.g., [2], §2.1).

So far we have defined some characteristics of the distribution P itself; now we will give definitions related to individual region predictors. The most natural class of region predictors is that of *permutationally invariant region predictors* Γ , for which $\Gamma_{1-\delta}(z_1, \dots, z_n, x)$ does not depend on the order of z_1, \dots, z_n (we know the examples are i.i.d., so knowing the order should not help).

The *calibration curve* of a region predictor Γ under P is the following function of the type $[0, 1] \rightarrow [0, 1]$:

$$\mathbf{C}(\delta) := \inf \left\{ \beta : \mathbb{P} \left\{ \limsup_{n \rightarrow \infty} \frac{\text{Err}_n(P^\infty, \Gamma_{1-\delta})}{n} \leq \beta \right\} = 1 \right\} \quad (2)$$

($\mathbb{P}(E)$ stands for the probability of event E). By the Hewitt–Savage zero-one law (see, e.g., [6], Theorem IV.1.3) in the case of permutationally invariant region predictors this definition will not change if “= 1” is replaced by “> 0” in (2). The *performance curve* of Γ under P is defined by

$$\mathbf{P}(\delta) := \inf \left\{ \beta : \mathbb{P} \left\{ \limsup_{n \rightarrow \infty} \frac{\text{Unc}_n(P^\infty, \Gamma_{1-\delta})}{n} \leq \beta \right\} = 1 \right\}; \quad (3)$$

this is again a function of the type $[0, 1] \rightarrow [0, 1]$. The Hewitt–Savage zero-one law again implies that for permutationally invariant Γ this will not change if “= 1” is replaced by “> 0”.

We will say that a region predictor Γ is *well-calibrated* under P if its calibration curve $\mathbf{C}(\delta)$ is below the diagonal: $\mathbf{C}(\delta) \leq \delta$ for any significance level δ . It is *asymptotically optimal* under P if its performance curve coincides with the success curve: $\mathbf{P}(\delta) = \mathbf{S}(\delta)$ for all δ .

Theorem 1 *Let P be a probability distribution in \mathbf{Z} with success curve \mathbf{S} . If a region predictor Γ is well-calibrated under P , its performance curve \mathbf{P} is above \mathbf{S} : for any δ , $\mathbf{P}(\delta) \geq \mathbf{S}(\delta)$. Moreover, for any significance level δ ,*

$$\liminf_{n \rightarrow \infty} \frac{\text{Unc}_n(P^\infty, \Gamma_{1-\delta})}{n} \geq \mathbf{S}(\delta) \quad a.s. \quad (4)$$

Let us now assume, for simplicity, that the distribution P is *regular*, in the sense that the predictability distribution function F is continuous. (The general case will be considered in §5 and will involve randomised region predictors.)

The main result of this paper (Theorem 2, strengthened in Theorem 2r) is that one can construct an asymptotically optimal TCM (which is well-calibrated automatically, by [10]). If, however, we know for sure that the true distribution is P it is very easy to construct a well-calibrated and asymptotically optimal region predictor. Fix a *choice function* $\hat{y} : \mathbf{X} \rightarrow \mathbf{Y}$ such that

$$\forall x \in \mathbf{X} : f(x) = P(\hat{y}(x) | x)$$

(to put it differently, $\hat{y}(x) \in \arg \max_y P(y | x)$). Define the *P-Bayesian* region predictor Γ by

$$\Gamma_{1-\delta}(z_1, \dots, z_n, x) := \begin{cases} \{\hat{y}(x)\} & \text{if } F(f(x)) \geq \mathbf{S}(\delta) \\ \mathbf{Y} & \text{otherwise,} \end{cases}$$

for all significance levels δ and data sequences $(z_1, \dots, z_n, x) \in \mathbf{Z}^n \times \mathbf{X}$, $n = 0, 1, \dots$. It can be shown that the *P*-Bayesian region predictor is well-calibrated and asymptotically optimal under P . (Our definition of the *P*-Bayesian region predictor is arbitrary in several respects; in principle, different choice functions can be used at different trials, the prediction can be arbitrary when $F(f(x)) = \mathbf{S}(\delta)$, and \mathbf{Y} can be replaced by any $E \subseteq \mathbf{Y}$ such that $P(E | x) := \sum_{y \in E} P(y | x) = 1$.)

The critical significance level (1) is an important characteristic of the probability distribution P generating the examples. If $\delta > \delta_0$, an optimal region predictor will always output certain predictions and, if forced to achieve the error rate δ , will sometimes have to output empty predictions. If, on the other hand, $\delta < \delta_0$, there will be uncertain predictions but no empty predictions. Figure 1 suggests that the critical significance level for the USPS data set is between 1% and 5%. This agrees with the observation that the critical significance level is just the error rate of the Bayesian point predictor (which is restricted to outputting Γ_n with $|\Gamma_n| = 1$ and minimises the expected number of errors) and the fact (reported in [7]) that the error rate achieved by humans on the USPS data set is 2.5%. Notice that in Figure 1 (left) the onset of empty predictions closely follows the point where all predictions become certain; see also Figures 4 and 5.

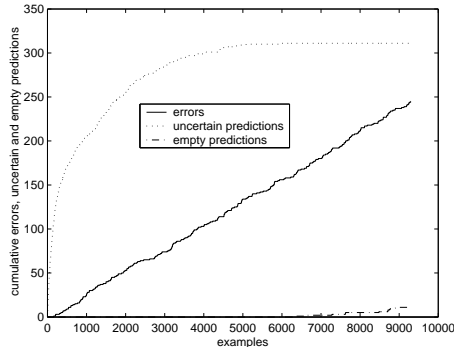


Figure 4: On-line performance of the Nearest Neighbour TCM on the USPS data set for the confidence level 97.5%. (This figure cannot be directly compared to the error rate of 2.5% for humans reported in [7], since this experiment has been carried out on the randomly permuted data set, whereas the test part of the USPS data set is known to be especially hard.)

3 Transductive Confidence Machine

The procedure at the end of §2 works well when P is known. If, however, P is only a convenient benchmark, the Bayesian region predictor can give very misleading results [4]. In the rest of this paper we will discuss how to ensure well-calibratedness under any distribution in \mathbf{Z} without losing the asymptotic performance of the Bayesian predictor if P happens to be the true distribution.

Transductive Confidence Machine (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 and \mathbb{R} is the set of all real numbers (equipped with the Borel σ -algebra), is called an *individual strangeness measure* if, for any $n = 1, 2, \dots$, each α_i in

$$A_n : (z_1, \dots, z_n) \mapsto (\alpha_1, \dots, \alpha_n) \tag{5}$$

is determined by z_i and the bag $\{z_1, \dots, z_n\}$. (The difference between the bag $\{z_1, \dots, z_n\}$ and the set $\{z_1, \dots, z_n\}$ is that the former can contain several copies of the same element.)

The *TCM associated with the individual strangeness measure A_n* is the following region predictor: $\Gamma_{1-\delta}(x_1, y_1, \dots, x_{n-1}, y_{n-1}, x_n)$ is defined to be

the set of all labels $y \in \mathbf{Y}$ such that

$$\frac{\#\{i = 1, \dots, n : \alpha_i \geq \alpha_n\}}{n} > \delta, \quad (6)$$

where

$$(\alpha_1, \dots, \alpha_n) := A_n((x_1, y_1), \dots, (x_{n-1}, y_{n-1}), (x_n, y)). \quad (7)$$

In general, a *TCM* is the TCM associated with some individual strangeness measure. It is shown in [10] that TCM is well-calibrated under any P (the technical report [10] also contains stronger assertions: for example, TCM is still well-calibrated, in a natural sense, when the confidence level $1 - \delta$ is allowed to depend on n).

4 Asymptotically Optimal TCM

If we suspect that the probability distribution in \mathbf{Z} generating the examples might be P , we can define the individual strangeness measure (5) by

$$\alpha_i := \begin{cases} 0 & \text{if } y_i = \hat{y}(x_i) \\ P(\hat{y}(x_i) | x_i) & \text{otherwise.} \end{cases} \quad (8)$$

The corresponding TCM will be called the *P-TCM* (cf. [9]). We say that a region predictor is *universally well-calibrated* if it is well-calibrated under any probability distribution P in \mathbf{Z} .

Theorem 2 *Let P be a regular probability distribution in \mathbf{Z} . The P -TCM is (a) universally well-calibrated and (b) asymptotically optimal under P .*

5 Randomised Region Predictors

In this section we will remove the assumption that the probability distribution P generating examples is regular. The price we will have to pay is that we will have to generalise the notion of region predictor in general, and TCM in particular, to allow using a generator of random numbers (as in [10]).

The generator that we consider generates a sequence τ_n , $n = 1, 2, \dots$, of uniformly distributed independent random numbers in the interval $[0, 1]$; τ_n will be used by Predictor at trial n of our basic protocol (see §1). Formally, a *randomised region predictor* Γ is a family, indexed by $n = 1, 2, \dots$ and

$\gamma \in [0, 1]$, of measurable functions $\Gamma_\gamma(z_1, \tau_1, \dots, z_{n-1}, \tau_{n-1}, x_n, \tau_n)$, where the $z_i \in \mathbf{Z}$, $i = 1, \dots, n-1$, are examples, $\tau_i \in [0, 1]$, $i = 1, \dots, n$, and $x_n \in \mathbf{X}$ is an object, which satisfies

$$\Gamma_{\gamma_1}(z_1, \tau_1, \dots, z_{n-1}, \tau_{n-1}, x_n, \tau_n) \subseteq \Gamma_{\gamma_2}(z_1, \tau_1, \dots, z_{n-1}, \tau_{n-1}, x_n, \tau_n)$$

whenever $\gamma_1 \leq \gamma_2$. The notation err_n , unc_n , etc., will be continued to be used in the randomised case as well; it should be remembered that these now depend on τ_1, τ_2, \dots . We can strengthen Theorem 1 as follows.

Theorem 1r. *Let P be a probability distribution in \mathbf{Z} with success curve \mathbf{S} . If a randomised region predictor Γ is well-calibrated under P , then for any significance level δ ,*

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

The ‘‘a.s.’’ in this theorem refers to the probability distribution $(P \times \mathbf{U})^\infty$ generating the sequence $z_1, \tau_1, z_2, \tau_2, \dots$, with \mathbf{U} standing for the uniform distribution in $[0, 1]$.

Next we introduce a randomised version of TCM. The *randomised Transductive Confidence Machine (rTCM)* associated with an individual strangeness measure A_n is the following randomised region predictor $\Gamma_{1-\delta}$: at any trial n and for any label $y \in \mathbf{Y}$,

1. if $\#\{i = 1, \dots, n : \alpha_i > \alpha_n\}/n > \delta$ (as before, the α s are defined by (7)), the label y is included in $\Gamma_{1-\delta}$;
2. if $\#\{i = 1, \dots, n : \alpha_i \geq \alpha_n\}/n \leq \delta$, y is not included in $\Gamma_{1-\delta}$;
3. otherwise, y is included in $\Gamma_{1-\delta}$ if

$$\tau_n < \frac{\#\{i = 1, \dots, n : \alpha_i \geq \alpha_n\} - n\delta}{\#\{i = 1, \dots, n : \alpha_i = \alpha_n\}}. \quad (9)$$

We say that a randomised region predictor Γ is *perfectly calibrated* if, under any probability distribution P , its calibration curve $\mathbf{C}(\delta)$ coincides with the diagonal: $\mathbf{C}(\delta) = \delta$ for any significance level δ . (Formally, this definition can also be given for deterministic predictors as well, but it would be impossible to satisfy in some cases.) The P -rTCM is defined to be the rTCM associated with the individual strangeness measure (8).

Theorem 2r. *Let P be a probability distribution in \mathbf{Z} . The P -rTCM is perfectly calibrated and, under P , asymptotically optimal.*

Notice that the rTCM makes at least as many errors as the TCM associated with the same individual strangeness measure. It is shown in [10] that rTCM's errors are independent and happen with probability δ at any confidence level $1 - \delta$. The difference between TCM and rTCM is typically negligible after the first several hundred trials; cf. the dotted line in Figure 2 (left).

6 Conclusion

In this paper we defined two desiderata for region predictors: being well-calibrated and asymptotic optimality. Being well-calibrated is the first priority: without it, the meaning of confidence levels is lost and it does not make sense to talk about optimality. If the probability distribution P generating individual examples is known, the Bayesian region predictor is well-calibrated and asymptotically optimal. But even in the situation where P is just a convenient guess that we are unwilling to take too seriously, there exists a region predictor (a TCM) which is universally well-calibrated and asymptotically optimal under P .

Acknowledgments

I am grateful to Alex Gammerman, Philip Dawid and anonymous referees for helpful suggestions. This work was partially supported by EPSRC (grant GR/R46670/01 “Complexity Approximation Principle and Predictive Complexity: Analysis and Applications”), BBSRC (grant 111/BIO14428 “Pattern Recognition Techniques for Gene and Promoter Identification and Classification in Plant Genomic Sequences”), and EU (grant IST-1999-10226 “EurEdit: The Development and Evaluation of New Methods for Editing and Imputation”).

References

- [1] Nicolas Bourbaki. *Eléments de mathématique, Livre IV, Fonctions d'une variable réelle (théorie élémentaire)*. Second edition. Hermann, Paris

- (1958).
- [2] Luc Devroye, László Györfi, and Gábor Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer, New York (1996).
 - [3] E. Lehmann. *Testing Statistical Hypotheses*. Wiley, New York (1959)
 - [4] Tom Melliush, Craig Saunders, Ilia Nouretdinov, and Vladimir Vovk. Comparing the Bayes and typicalness frameworks. In: L. De Raedt and P. Flash (eds.), *Machine Learning: ECML 2001. Proceedings of the Twelfth European Conference on Machine Learning*. Lecture Notes in Artificial Intelligence, Vol. 2167, Springer (2001) 360–371. Full version published as Technical Report CLRC-TR-01-05, Computer Learning Research Centre, Royal Holloway, University of London, <http://www.clrc.rhul.ac.uk>.
 - [5] Craig Saunders, Alex Gammerman, and Vladimir Vovk. Transduction with confidence and credibility. In: *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence* (1999) 722–726.
 - [6] Albert N. Shiryaev. *Probability*. Second edition. Springer, New York (1996).
 - [7] Vladimir N. Vapnik. *Statistical Learning Theory*. Wiley, New York (1998).
 - [8] Vladimir Vovk, Alex Gammerman, and Craig Saunders. Machine-learning applications of algorithmic randomness. In: *Proceedings of the Sixteenth International Conference on Machine Learning* (1999) 444–453.
 - [9] Vladimir Vovk. De-Bayesing Bayesian prediction algorithms. Manuscript (June 2000).
 - [10] Vladimir Vovk. On-line Confidence Machines are well-calibrated. On-line Compression Modelling project, Working Paper #1. In: *Proceedings of FOCS'2002*. Full version published as Technical Report CLRC-TR-02-01, Computer Learning Research Centre, Royal Holloway, University of London, <http://www.clrc.rhul.ac.uk> (April 2002). Additional information can be found at <http://vovk.net/kp>.

Appendix: Proofs

First we establish some simple properties of the predictability distribution function and success curve.

Lemma 1 *The predictability distribution function F satisfies the following properties:*

1. $F(\epsilon) = 0$ for some $\epsilon > 0$ and $F(1) = 1$;
2. F is non-decreasing;
3. F is continuous on the right.

If a function $F : [0, 1] \rightarrow [0, 1]$ 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 F is the predictability distribution function.

Proof Properties 1 (cf. the caption to Figure 3), 2, and 3 are obvious (and the last two are well-known properties of all distribution functions). The fact that these three properties characterise predictability distribution functions easily follows from the fact that the last two properties plus $F(-\infty) = 0$ and $F(\infty) = 1$ characterise distribution functions (see, e.g., [6], Theorem II.3.1). ■

We will use the notations g'_{left} and g'_{right} for the left and right derivatives, respectively, of a function g .

Lemma 2 *The success curve $\mathbf{S} : [0, 1] \rightarrow [0, 1]$ always satisfies these properties:*

1. \mathbf{S} is convex.
2. *There is a point $\delta_0 \in [0, 1]$ (the critical significance level) such that $\mathbf{S}(\delta) = 0$ for $\delta \geq \delta_0$ and $\mathbf{S}'_{\text{left}}(\delta_0) < -1$; therefore, $\mathbf{S}'_{\text{left}} < -1$ and $\mathbf{S}'_{\text{right}} < -1$ to the left of δ_0 , and the function \mathbf{S} is decreasing before it hits the δ -axis at δ_0 .*
3. \mathbf{S} is continuous at $\delta = 0$; therefore, it is continuous everywhere in $[0, 1]$.

If a function $\mathbf{S} : [0, 1] \rightarrow [0, 1]$ 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{S} is the success curve.

Proof For the basic properties of convex functions and their left and right derivatives, see, e.g., [1], §I.4. The statement of the lemma follows from the fact that the success curve \mathbf{S} can be obtained from the predictability distribution function F using these steps (labelling the horizontal and vertical axes as x and y respectively):

1. Invert F : $F_1 := F^{-1}$.
2. Flip F_1 around the line $x = 0.5$ and then around the line $y = 0.5$:
 $F_2(x) := 1 - F_1(1 - x)$.
3. Integrate F_2 : $F_3(x) := \int_0^x F_2(t)dt$.
4. Invert F_3 : $F_4 := F_3^{-1}$.
5. Flip F_4 around the line $y = 0.5$: $F_5 := 1 - F_4$.

It can be shown that $\mathbf{S} = F_5$, no matter which of the several natural definitions of the operation $g \mapsto g^{-1}$ is used; for concreteness, we can define $g^{-1}(y) := \sup\{x : g(x) \leq y\}$ for non-decreasing g (so that g^{-1} is continuous on the right). ■

Visually the empirical performance curve in Figure 2 seems to satisfy the properties listed in Lemma 2 for significance levels that are not too large or too small (approximately in the range 0.1%–5%); for an even better agreement, see Figure 5.

A natural idea is to reverse the process of transforming F into \mathbf{S} and try to obtain an estimate of the predictability distribution function F from an empirical performance curve. Figure 6 shows the result of such an attempt. Such pictures, however, should not be taken too seriously, since the differentiation operation needed in finding F is known to be unstable (see, e.g., [7], §1.12).

Proof of Theorems 1 and 1r.

Let us check first that (4) indeed implies $\mathbf{P}(\delta) \geq \mathbf{S}(\delta)$. Since probability measures are σ -additive, (3) implies

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

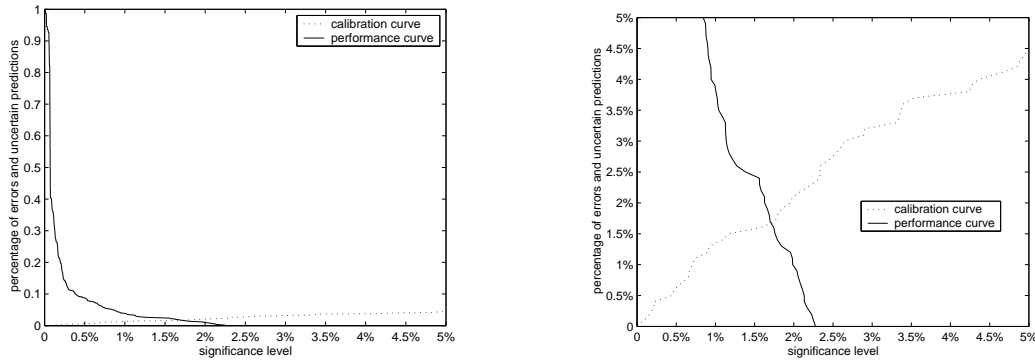


Figure 5: **Left:** Picture analogous to Figure 2 (right) for the last one thousand examples. Notice a different behaviour of the empirical performance curve as it approaches the horizontal axis as compared with Figure 2 (right). The unexpected behaviour of the empirical performance curve as it approaches the vertical axis may be explained (at least partially) by the “granularity” of TCM: for example, the “realised p-value” given by the left-hand side of (6) can never be less than $1/9298 > 0.01\%$; this behaviour may become more regular for randomised TCM. **Right:** The bottom part of the picture on the left stretched vertically. Notice that the slope of the empirical performance curve is at least 1 in absolute value before it hits the horizontal axis; this agrees with Lemma 2 on p. 12. This figure suggests that, if the 1-NN TCM were an optimal region predictor, the critical significance level for the USPS data set would be close to 2.3.

and so we obtain from (4):

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

almost surely; since the two extreme terms are deterministic, we have $\mathbf{P}(\delta) \geq \mathbf{S}(\delta)$.

We start the actual proof with alternative definitions of calibration and performance curves. Complement the protocol of §1 in which Nature plays P^∞ and Predictor plays $\Gamma_{1-\delta}$ with the following variables:

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

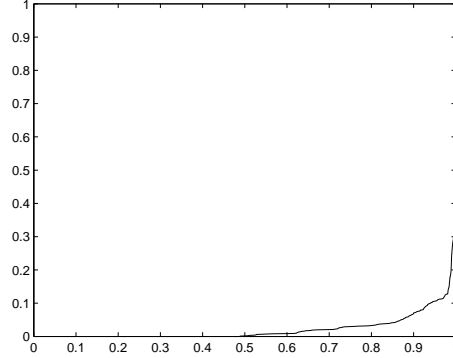


Figure 6: An attempt to reverse engineer the predictability distribution function of the hand-written digits in the USPS data set. This picture was obtained from the solid line in Figure 5 (left) by reversing the list in the proof of Lemma 2.

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

(we are not always consistent in the order of arguments of the function $\Gamma_{1-\delta}$). The *prequential calibration curve* of Γ under P is defined by

$$\overline{\mathbf{C}}(\delta) := \inf \left\{ \beta : \mathbb{P} \left\{ \limsup_{n \rightarrow \infty} \frac{\overline{\text{Err}}_n(P^\infty, \Gamma_{1-\delta})}{n} \leq \beta \right\} = 1 \right\}$$

and the *prequential performance curve* of Γ under P by

$$\overline{\mathbf{P}}(\delta) := \inf \left\{ \beta : \mathbb{P} \left\{ \limsup_{n \rightarrow \infty} \frac{\overline{\text{Unc}}_n(P^\infty, \Gamma_{1-\delta})}{n} \leq \beta \right\} = 1 \right\},$$

where \mathbb{P} refers to the probability distribution $(P \times \mathbf{U})^\infty$ over the examples z_1, z_2, \dots and random numbers τ_1, τ_2, \dots . By the martingale strong law of large numbers the prequential versions of the calibration and performance curves coincide with the original versions: 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 z_1, \dots, z_n and τ_1, \dots, τ_n , we have

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

and

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

(see, e.g., [6], Theorem VII.5.4). It is also clear that we can replace Unc_n by $\overline{\text{Unc}}_n$ in (4).

Without loss of generality we can assume that Nature's move Γ_n at trial n is either $\{\hat{y}(x_n)\}$ or the whole label space \mathbf{Y} . Furthermore, we can assume that

$$\overline{\text{unc}}_n = \mathbf{S}(\overline{\text{err}}_n) \tag{10}$$

at every trial, since the best way to spend the allowance of $\overline{\text{err}}_n$ is to be certain on objects x with the largest (topmost in Figure 3) representations $F(f(x))$. (For a formal argument, see the end of this proof.) Using the fact that the success curve \mathbf{S} is convex, non-increasing, and continuous (see Lemma 2), we obtain

$$\frac{\overline{\text{Unc}}_n}{n} = \frac{1}{n} \sum_{i=1}^n \overline{\text{unc}}_i = \frac{1}{n} \sum_{i=1}^n \mathbf{S}(\overline{\text{err}}_i) \geq \mathbf{S}\left(\frac{1}{n} \sum_{i=1}^n \overline{\text{err}}_i\right) = \mathbf{S}\left(\frac{\overline{\text{Err}}_n}{n}\right) \geq \mathbf{S}(\delta) - \epsilon,$$

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

It remains to prove formally that $\overline{\text{unc}}_n \geq \mathbf{S}(\overline{\text{err}}_n)$ (which is the part of (10) that we actually used). Let us fix $1 - \delta$ and $x_1, \tau_1, y_1, \dots, x_{n-1}, \tau_{n-1}, y_{n-1}$; we will write

$$\Gamma(x, \tau) := \Gamma_{1-\delta}(x_1, \tau_1, y_1, \dots, x_{n-1}, \tau_{n-1}, y_{n-1}, x, \tau),$$

omitting the fixed arguments. Without loss of generality we are assuming that either $\Gamma(x, \tau) = \hat{y}(x)$ or $\Gamma(x, \tau) = \mathbf{Y}$. Set

$$p(x) := \mathbf{U} \{ \tau : \Gamma(x, \tau) = \{\hat{y}(x)\} \}, \quad \delta := \overline{\text{err}}_n.$$

Our goal is to show that $\overline{\text{unc}}_n \geq \mathbf{S}(\delta)$; without loss of generality we assume $0 < \delta < \delta_0$. To put it differently, we are required to show that the value of the optimisation problem

$$\int_{\mathbf{X}} p(x) P(dx) \rightarrow \max \tag{11}$$

subject to the constraint

$$\int_{\mathbf{X}} (1 - f(x)) p(x) P(dx) = \delta$$

is $1 - \mathbf{S}(\delta)$ at best. By the Neyman–Pearson lemma (see, e.g., [3]) there exist constants $c > 0$ and $d \in [0, 1]$ such that

$$p(x) = \begin{cases} 1 & \text{if } f(x) > c \\ d & \text{if } f(x) = c \\ 0 & \text{if } f(x) < c . \end{cases} \quad (12)$$

The constants c and d are defined (c uniquely and d uniquely unless the probability of $f(x) = c$ is zero) from the condition

$$\int_{x:f(x)>c} (1 - f(x))P(dx) + d \int_{x:f(x)=c} (1 - c)P(dx) = \delta ,$$

which is equivalent, by Fubini’s theorem (applied to the indicator function of the subgraph of F ; see Figure 3, left), to

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

where $F(c-)$ is defined as $\lim_{\beta \uparrow c} F(\beta)$. From this it is easy to obtain that the value of the optimal problem (11) is indeed $1 - \mathbf{S}(\delta)$: using the notation $p_d(x)$ for the right-hand side of (12), we have

$$\begin{aligned} \int_{\mathbf{X}} p_d(x)P(dx) &= d \int p_1(x)P(dx) + (1 - d) \int p_0(x)P(dx) \\ &= dP\{x : f(x) \geq c\} + (1 - d)P\{x : f(x) > c\} \\ &= d(1 - F(c-)) + (1 - d)(1 - F(c)) \\ &= 1 - F(c) + d(F(c) - F(c-)) \\ &= 1 - \mathbf{S}(\delta) . \end{aligned}$$

This completes the proof of Theorem 1r; since Theorem 1 is a special case, it is also proved.

Proof of Theorems 2 and 2r.

The fact that every rTCM Γ is perfectly calibrated is proved in [10], so we are only required to show that Γ is asymptotically optimal under P . Fix a confidence level $1 - \delta$; we will show that

$$\limsup_{n \rightarrow \infty} \frac{\overline{\text{Unc}}_n(P^\infty, \Gamma_{1-\delta})}{n} \leq \mathbf{S}(\delta) \quad (13)$$

almost surely (the underlying probability distribution \mathbb{P} being the product $(P \times \mathbf{U})^\infty$). Without loss of generality we assume $\mathbf{S}(\delta) < 1$ ((13) holds trivially when $\mathbf{S}(\delta) = 1$). Set

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

The case $c = 1$ is simple: it means that $P\{x : f(x) < 1\} \leq \mathbf{S}(\delta)$; since, almost surely, $\text{unc}_n = 0$ at trials where $f(x_n) = 1$, by Borel's strong law of large numbers we immediately obtain (13). Therefore, we assume $c < 1$ in the rest of the proof.

First we consider the case $F(c) = \mathbf{S}(\delta)$ (this will be sufficient to prove Theorem 2, since $F(c) = \mathbf{S}(\delta)$ is implied by $F(c) = F(c-)$ and the rTCM constructed in this part of the proof will be deterministic). Notice that $F(c + \epsilon) > F(c)$ for any $0 < \epsilon \leq 1 - c$ (we are assuming $\epsilon \leq 1 - c$ so that $F(c + \epsilon)$ is defined). We will prove that, for any $0 < \epsilon \leq 1 - c$ and from some n on,

$$\mathbb{P}(\text{unc}_n \mid \mathcal{F}_{n-1}) \leq F(c + \epsilon) \quad \text{a.s.} \quad (14)$$

(we are using the same notation for an event and for its indicator function). This will imply

$$\limsup_{n \rightarrow \infty} \frac{\overline{\text{Unc}_n}}{n} \leq F(c + \epsilon)$$

almost surely; since $\lim_{\epsilon \downarrow 0} F(c + \epsilon) = \mathbf{S}(\delta)$, this will prove (13).

Fix $0 < \epsilon \leq 1 - c$; without loss of generality assume that F is continuous at $c + \epsilon$. Let us prove (14), assuming n is large enough. Suppose the examples observed before trial n are $(z_1, \dots, z_{n-1}) = ((x_1, y_1), \dots, (x_{n-1}, y_{n-1}))$. Let us say that an example $(x, y) \in \mathbf{Z}$ is *wrongly classified* if $y \neq \hat{y}(x)$. Remember that according to the individual strangeness measure (8) the strange elements in every bag of examples are those that are wrongly classified, and the more predictable they are the stranger. Notice that the prediction for the new object x_n will be certain if (a) the new object x_n is to the right of $\beta = c + \epsilon$ in Figure 3 (left), in the sense $f(x_n) \geq c + \epsilon$, and (b) the number of the wrongly classified objects x_i , $i = 1, \dots, n - 1$, to the right of $\beta = c + \epsilon$ is less than $n\delta - 10$. The probability (conditional on \mathcal{F}_{n-1}) of (a) is $1 - F(c + \epsilon)$, so to prove (14) it is sufficient to show that the event (b) (remember that it is measurable w.r. to \mathcal{F}_{n-1}) happens from some n on almost surely. The probability that an object is wrongly classified and to the right of $\beta = c + \epsilon$ is

$$b := \int_0^1 (F(\beta) - F(c + \epsilon))^+ d\beta < \delta .$$

By Hoeffding's inequality (see, e.g., [2], Theorem 8.1) the probability that the event (b) will fail to happen is bounded from above by

$$e^{-2(n\delta-10-(n-1)b)^2/(n-1)} \leq e^{-\kappa n}, \quad (15)$$

for some positive constant κ and from some n on. Since $\sum_n e^{-\kappa n} < \infty$, the Borel–Cantelli lemma implies that (b) will almost surely happen from some n on. This completes the proof in the case $F(c) = \mathbf{S}(\delta)$.

Now we consider the case $F(c) > \mathbf{S}(\delta)$ (which is the only remaining possibility). In the remaining part of the proof it will be important that we consider rTCM rather than TCM.

Let $\epsilon > 0$ satisfy $\epsilon < F(c) - \mathbf{S}(\delta)$. We will prove that, from some n on,

$$\mathbb{P}(\text{unc}_n | \mathcal{F}_{n-1}) \leq \mathbf{S}(\delta) + \epsilon \quad \text{a.s.} \quad (16)$$

This will imply

$$\limsup_{n \rightarrow \infty} \frac{\overline{\text{Unc}}_n}{n} \leq \mathbf{S}(\delta) + \epsilon$$

almost surely, and so prove (13).

We say that an object and random number $(x, \tau) \in \mathbf{X} \times [0, 1]$ (such a pair will be called an *extended object*) is *above the line* $\mathbf{S}(\delta) + \epsilon$ (cf. Figure 3, left) if either $f(x) > c$ or

$$f(x) = c \ \& \ \tau \geq \frac{\mathbf{S}(\delta) + \epsilon - F(c-)}{F(c) - F(c-)}$$

(this definition corresponds to representing each extended object (x, τ) by the point

$$(f(x), \tau F(f(x)) + (1 - \tau)F(f(x)-))$$

in Figure 3, left).

Let us prove (16), assuming n is large enough. Suppose the extended objects observed before trial n are $(x_1, \tau_1, \dots, x_{n-1}, \tau_{n-1})$. Now the prediction for the new object x_n will be certain if (a) the new extended object (x_n, τ_n) is above $\mathbf{S}(\delta) + \epsilon$, and (b) the number of the wrongly classified extended objects (x_i, τ_i) , $i = 1, \dots, n-1$, above $\mathbf{S}(\delta) + \epsilon$ is less than $n\delta - 10$. (We say that (x, τ) is wrongly classified if x is.) The probability (conditional on \mathcal{F}_{n-1}) of (a) is $1 - \mathbf{S}(\delta) - \epsilon$, so to prove (16) it is sufficient to show that the

event (b) happens from some n on almost surely. The probability that an extended object is wrongly classified and above $\mathbf{S}(\delta) + \epsilon$ is

$$b := \int_0^1 (F(\beta) - \mathbf{S}(\delta) - \epsilon)^+ d\beta < \delta .$$

The proof is completed literally as before: apply Hoeffding's inequality to obtain upper bound (15) and then apply the Borel–Cantelli lemma.

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 Nourtdinov and Alex Gammerman, March 2003.
5. *Testing exchangeability on-line*, by Vladimir Vovk, Ilia Nourtdinov and Alex Gammerman, February 2003.
6. *Criterion of calibration for Transductive Confidence Machine with limited feedback*, by Ilia Nourtdinov 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.