

Cross-validation for binary classification by real-valued functions: theoretical analysis

Martin Anthony
Department of Mathematics
London School of Economics
Houghton Street
London WC2A 2AE

Sean B. Holden
Department of Computer Science
University College London
Gower Street
London WC1E 6BT

DRAFT

Abstract

This paper concerns the use of real-valued functions for binary classification problems. Previous work in this area has concentrated on using as an error estimate the ‘resubstitution’ error (that is, the empirical error of a classifier on the training sample) or its derivatives. However, in practice, cross-validation and related techniques are more popular. Here, we devise new holdout and cross-validation estimators for the case where real-valued functions are used as classifiers, and we analyse theoretically the accuracy of these.

Contents

1	Introduction	3
2	Error estimation	4
3	Previous related work	5
3.1	PAC learning in the realisable case	6

3.2	Empirical errors	8
3.3	The holdout estimate	9
3.4	The cross-validation estimate	9
3.5	Previous bounds	9
4	Using real-valued function classes	11
4.1	Measuring error with a margin	11
4.2	Measuring cross-validation error	12
5	Results for real-valued cross-validation models	13
5.1	Covering numbers	13
5.2	Main probability theorems	14
5.3	Discussion	16
6	New estimators	17
6.1	Modified estimators	17
6.2	Performance bounds	18
6.3	Application to neural networks	19
7	General ‘dimension-based’ bounds	20
8	Proofs	23
8.1	Relating cross-validation estimate to holdout estimate	23
8.2	Symmetrization	24
8.3	Using a group action	25

8.4	Using covers	26
8.5	Proof of the main theorems	28
9	Conclusions and further work	30

1 Introduction

Recently, within probabilistic models of machine learning, attention has focused on the use of real-valued functions for binary classification (as in [22, 8, 9, 4], for instance). It has been shown that in many cases, one can obtain fairly accurate estimates of a classifier’s error when the classifier is a real-valued functions which achieves the correct classification of a training sample, with a large ‘margin’ (a notion to be made precise in what follows). Not only have the results obtained for such classification models often improved upon the analogous results for classification by binary-valued classes: in some cases, error estimates can be given in the new model when none can possibly be given in the binary classification model. In this paper we take a further step in the analysis of classification by real-valued classes. Previous analyses have concentrated on using as an error estimate the ‘resubstitution’ error (that is, the empirical error of a classifier on the training sample) or its derivatives. However, in practice, cross-validation and related techniques are more popular. In this paper, we devise a new technique for cross-validation in the case where real-valued functions are used as classifiers. We analyse the performance of this technique probabilistically. The results we obtain are what Kearns and Ron [18] (in the binary-valued context) have described as ‘sanity-check’ bounds. That is, the bounds are not better than those one has for the corresponding resubstitution error estimate. Nonetheless, we believe that such sanity-check bounds are worthwhile. They show that our new cross-validation technique works. The error bounds might be looser than one would like for real, practical algorithms, but that is a consequence of the very general approach we take. The possibility is left open of our technique out-performing the resubstitution technique in practice.

We start, in Section 2 with a general discussion of the error estimation problem. In Section 3 we present previously obtained results to set the present work in context. In this section, we provide the basic definitions of the standard holdout and cross-validation estimators, and we describe the PAC probabilistic learning model. The next section 4 then discusses the bases of the new work: the use of real-valued functions for binary classification, and the importance of a ‘margin’ in the classification. We tentatively suggest a modified way of measuring holdout and cross-validation errors, taking into account the classification margin. Section 5 presents the probability results which are at the centre of our analysis (the proofs of which are deferred until Section 8).

Section 6, building on the results of Section 5 and the suggestions made in Section 4, formalises the definitions of two new error estimators applicable when using real-valued functions for classification purposes and discusses their accuracy. Applications to certain types of neural network are given by way of illustration. In Section 7 we obtain general error bounds in terms of the ‘fat-shattering dimension’, a parameter describing the complexity of a function class. Section 8 contains the proofs of the probability results of Section 5. The concluding Section 9 suggests some open problems.

2 Error estimation

In this section we describe the standard classification and cross-validation frameworks and present some previous results on the error estimation abilities of these techniques.

We start by describing a standard form of the ‘PAC’ model of learning [24, 10, 14]. This concerns two-class pattern classification problems, in which we wish to classify objects $x \in X$ as belonging to class -1 or class 1 . (The set X will be called the set of *examples*.) The appearance of examples and their class membership labels is governed by an arbitrary distribution P on $Z = X \times \{-1, 1\}$. During ‘learning’, we are presented with a *training sample* \mathbf{z} of n *labelled examples* drawn independently at random according to P ,

$$\mathbf{z} = ((x_1, y_1), \dots, (x_n, y_n)).$$

A learning algorithm L takes such samples \mathbf{z} and uses them to select a hypothesis $h : X \rightarrow \{-1, 1\}$ from a hypothesis space \mathcal{H} , which may be used to classify further objects drawn independently at random according to P . The (true) *error* $\text{er}_P(h)$ of $h = L(\mathbf{z})$ is defined to be

$$\text{er}_P(h) = P\{(x, y) \in Z : h(x) \neq y\}$$

and the *empirical error* $\text{er}_{\mathbf{z}}(h)$ is defined to be

$$\text{er}_{\mathbf{z}}(h) = \frac{1}{n} |\{i : h(x_i) \neq y_i\}|.$$

The empirical error is often referred to as the ‘error on the training sample’ or, in statistics, the ‘resubstitution estimate’, and can be regarded as an estimate of the value of the actual error $\text{er}_P(h)$. It is natural to ask how close it is to the actual error, and a significant body of research exists which bounds the probability that the resubstitution estimate is a ‘bad’ estimate of the actual error in an appropriately defined sense. A typical result is the following, due to Vapnik [25] (see also [10]). (Like all such results, this requires some fairly reasonable measurability conditions on the hypothesis space \mathcal{H} . These are not very stringent, and we shall assume without further comment that they hold.)

Theorem 1 For parameters $\epsilon, \gamma \in (0, 1]$, and for any distribution P , we have

$$P^n \{ \mathbf{z} \in Z^n : \exists h \in \mathcal{H} \text{ for which } \text{er}_{\mathbf{z}}(h) \leq (1 - \gamma)\epsilon \text{ and } \text{er}_P(h) > \epsilon \} < 8\Pi_{\mathcal{H}}(2n) \exp(-\gamma^2\epsilon n/4)$$

where $\Pi_{\mathcal{H}}()$ denotes the growth function which we define below.

Theorem 1 tells us that the probability of obtaining a training sequence such that L can choose a hypothesis from \mathcal{H} that simultaneously has empirical error no greater than $(1 - \gamma)\epsilon$ and true error greater than ϵ is bounded by the stated expression. If \mathcal{H} has finite *Vapnik-Chervonenkis (VC) dimension* (see for example [10, 3]) then the upper bound can be made arbitrarily small by choosing a large enough n (which does not depend on P); this will be explained below.

A great deal of research in recent years has concentrated on extending results of this kind to deal with multiple class problems (see for example Anthony and Shawe-Taylor [7]), and real-valued functions. (See for example Haussler [14], Pollard [20], Anthony [2], among others.) However there is another potential direction for extending such results which to date has received much less emphasis, although it is no less interesting. Cheng and Titterington [11] have raised the interesting question of what happens to the result stated in Theorem 1 if a cross-validation estimate is used in place of a resubstitution estimate. This problem has recently been considered by Holden [15] and Kearns and Ron [18]. There are different types of cross-validation estimate. These operate by splitting the sample \mathbf{z} into m sections. A section is removed and the remaining examples used in conjunction with L to select a hypothesis h from \mathcal{H} . The removed section is then used to estimate the error of h , using the obvious technique of finding the fraction of the examples in the removed section that h classifies incorrectly (that is, the empirical error on the removed section). In the standard cross-validation techniques, this is repeated for each section and the average of the m individual estimates is then used as a final estimate of the true error that will be obtained if L is applied on the complete training sequence.

It is now well-known (see for example Toussaint [23]) that the resubstitution estimate is not a good estimator of true error in practice due to its considerable optimistic bias. (This simply means that it tends to underestimate the true error.) Clearly therefore the study of error estimates other than the resubstitution estimate is of significant importance if we wish to obtain results applicable to practical learning applications.

3 Previous related work

We now describe some previous work of Holden [15] which derives results like Theorem 1 for cross-validation error estimates in the ‘realisable’ case. The new results of this paper

provide similar results in a model of cross-validation where real-valued functions are used for classification.

3.1 PAC learning in the realisable case

The results in [15] apply to the *realisable* learning case, in which (with probability one) only one label occurs with each example. Explicitly, the correct classification for each example is assigned by a *target concept* $c : X \rightarrow \{-1, 1\}$, which we assume belongs to our *hypothesis class* \mathcal{H} . In this case, the training sample $\mathbf{z}(\mathbf{x}, c) \in (X \times \{-1, 1\})^n$ of length n is obtained by drawing n inputs independently at random according to an arbitrary distribution μ on X in order to obtain a sequence $\mathbf{x} = (x_1, \dots, x_n) \in X^n$, and then forming the training sample

$$\mathbf{z} = \mathbf{z}(\mathbf{x}, c) = ((x_1, c(x_1)), \dots, (x_n, c(x_n))).$$

We denote by $\mathbf{z}_{\mathcal{H}}$ the set of all possible training sequences,

$$\mathbf{z}_{\mathcal{H}} = \{\mathbf{z}(\mathbf{x}, c) : c \in \mathcal{H}, \mathbf{x} \in X^n, n \geq 1\}.$$

In general, we will abbreviate $\mathbf{z}(\mathbf{x}, c)$ to \mathbf{z} when the underlying target concept and input sequence are clear from the context. The length of the training sequence is denoted by n throughout this paper. It should be noted that the realisable framework just described is a special case of the framework in which the labelled examples are chosen according to a joint probability distribution P on $Z = X \times \{-1, 1\}$; see [10, 7]. In learning from examples we attempt to use \mathbf{x} to select a hypothesis $h : X \rightarrow \{-1, 1\}$ from \mathcal{H} . A *learning algorithm* or *learner* L is a function $L : \mathbf{z}_{\mathcal{H}} \rightarrow \mathcal{H}$ that accomplishes this. A learning algorithm is *consistent* if for any $\mathbf{z} \in \mathbf{z}_{\mathcal{H}}$, L produces a hypothesis h for which $h(x_i) = c(x_i)$ for $1 \leq i \leq n$. Given any training sequence $\mathbf{z} \in \mathbf{z}_{\mathcal{H}}$ of length n we denote by $\mathcal{H}[\mathbf{z}]$ the set of all hypotheses in \mathcal{H} that are consistent with \mathbf{z} ,

$$\mathcal{H}[\mathbf{z}] = \{h \in \mathcal{H} : h(x_i) = c(x_i) \text{ for } 1 \leq i \leq n\}.$$

Finally, the *true error* (usually abbreviated to the *error*) of any hypothesis $h \in \mathcal{H}$ is denoted $\text{er}_{\mu}(h, c)$, and is defined as the probability that h disagrees with c on an input chosen at random according to μ :

$$\text{er}_{\mu}(h, c) = \mu\{x \in X : h(x) \neq c(x)\}.$$

We abbreviate $\text{er}_{\mu}(h, c)$ to $\text{er}_{\mu}(h)$ when the target concept is clear from the context.

In order to describe the results of Holden [15] and to move towards the results of this paper, we now define the *growth function* and *Vapnik-Chervonenkis (VC) dimension*. Let $S = \{x_1, \dots, x_i\}$ be a subset of X with $|S| = i$ and let F be a class of functions

$f : X \rightarrow \{-1, 1\}$. With each $f \in F$ we associate the set $A_f = \{x \in X : f(x) = 1\}$. We define the set $\Pi_F(S)$ as

$$\Pi_F(S) = \{S \cap A_f : f \in F\}$$

and we define the growth function $\Pi_F(j)$ by

$$\Pi_F(j) = \max\{|\Pi_F(S)| : S \subseteq X \text{ and } |S| = j\}.$$

The VC dimension of F is denoted $\text{VCdim}(F)$ and defined as

$$\text{VCdim}(F) = \max\{j : \Pi_F(j) = 2^j\}.$$

By ‘Sauer’s lemma’ (see [21, 10, 3]) we know that if $d = \text{VCdim}(F)$ is finite and $j \geq d \geq 1$, then

$$\Pi_F(j) < \left(\frac{ej}{d}\right)^d.$$

The following theorem is a standard result in the concept learning framework.

Theorem 2 ([10]; see also [3], **proposition 8.2.3**) *For any distribution μ , any $c \in \mathcal{H}$, any ϵ , and any $n \geq 8/\epsilon$,*

$$P^n\{\mathbf{x} \in X^n : \exists h \in \mathcal{H}[\mathbf{z}] \text{ such that } \text{er}_\mu(h) \geq \epsilon\} < 2\Pi_{\mathcal{H}}(2n)2^{-\epsilon n/2}.$$

Equivalently, for any n and $\delta \in (0, 1)$, and any distribution μ , with probability at least $1 - \delta$, a sample \mathbf{x} of length n satisfies: if $h \in H$ and $\text{er}_{\mathbf{x}}(h) = 0$, then

$$\text{er}_\mu(h) < \epsilon(n, \delta) = \frac{8}{n} \left(\log \left(\frac{2}{\delta} \right) \log(\Pi_H(2n)) \right),$$

where \log denotes binary logarithm.

A consequence of this theorem and Sauer’s Lemma is that if \mathcal{H} has finite VC dimension, then the estimation error $\epsilon(n, \delta)$ for a consistent algorithm can be made arbitrarily small by making n large enough.

We now take a more general perspective. Assume that we have an *error estimator* which takes the training sequence \mathbf{z} and learner L and uses them to make an estimate of $\text{er}_\mu(L(\mathbf{z}))$. We will denote the value of the estimate by $\text{est}(L, \mathbf{z})$. Holden [15] obtains upper bounds on the quantity

$$P^n\{\mathbf{x} \in X^n : |\text{est}(L, \mathbf{z}) - \text{er}_\mu(L(\mathbf{z}))| > \epsilon\}.$$

when the error estimators arise from cross-validation. Such bounds can, as we shall see, easily be converted into high-confidence bounds on the estimation error. We now describe these estimates in more detail.

FIGURE

Figure 1: Any sequence $\mathbf{a} \in A^n$ can be divided into m contiguous sections of equal size called *folds*. It is assumed that $m \geq 2$ and m divides n . Each fold contains $n' = (n/m)$ elements of \mathbf{a} .

3.2 Empirical errors

A slight abuse of set-theoretic notation will be used in describing and dealing with the error estimation techniques of interest. The following definitions are most easily explained diagrammatically (figure 1). Given any sequence $\mathbf{a} = (a_1, \dots, a_n) \in A^n$ for some set A we will speak of it as being divided into m *folds*, where m is an integer, $m \geq 2$, and m divides n . The folds form contiguous sections of equal size, each fold containing $n' = (n/m)$ elements of \mathbf{a} . We denote by \mathbf{a}_i the i th fold,

$$\mathbf{a}_i = (a_{(i-1)n'+1}, \dots, a_{in'})$$

and we denote by $\mathbf{a} \setminus \mathbf{a}_i$ the sequence obtained by removing the i th fold from \mathbf{a} ,

$$\mathbf{a} \setminus \mathbf{a}_i = (a_1, \dots, a_{(i-1)n'}, a_{in'+1}, \dots, a_n).$$

Given any two sequences $\mathbf{a} \in A^{n_1}$ and $\mathbf{b} \in A^{n_2}$ we denote by \mathbf{ab} the sequence

$$\mathbf{ab} = (a_1, \dots, a_{n_1}, b_1, \dots, b_{n_2}) \in A^{n_1+n_2}.$$

Clearly when \mathbf{a} is divided into m folds we have $\mathbf{a} = \mathbf{a}_1 \mathbf{a}_2 \dots \mathbf{a}_m$. If $\mathbf{x} \in X^n$ is any sequence obtained by drawing n inputs independently at random according to μ , then for a hypothesis $h \in \mathcal{H}$ and a target concept $c \in \mathcal{H}$ we define the *empirical error* $\text{er}_{\mathbf{x}}(h, c)$ of h as

$$\text{er}_{\mathbf{x}}(h, c) = \frac{1}{n} \sum_{i=1}^n I_{\{h(x_i) \neq c(x_i)\}}$$

where I denotes the indicator function. If \mathbf{x} is divided into m folds, then for m hypotheses h_1, \dots, h_m each of which is a member of \mathcal{H} , and a target concept $c \in \mathcal{H}$, we define the *m -fold empirical error* $\text{er}_{\mathbf{x}}(h_1, \dots, h_m, c)$ to be

$$\text{er}_{\mathbf{x}}(h_1, \dots, h_m, c) = \frac{1}{m} \sum_{i=1}^m \text{er}_{\mathbf{x}_i}(h_i, c) = \frac{1}{n} \sum_{i=1}^m \sum_{j=1}^{n'} I_{\{h_i(x_{(i-1)n'+j}) \neq c(x_{(i-1)n'+j})\}}$$

In general, we abbreviate $\text{er}_{\mathbf{x}}(h, c)$ and $\text{er}_{\mathbf{x}}(h_1, \dots, h_m, c)$ to $\text{er}_{\mathbf{x}}(h)$ and $\text{er}_{\mathbf{x}}(h_1, \dots, h_m)$ respectively when the target concept of interest is clear from the context.

3.3 The holdout estimate

Of particular importance in the theory developed in [15], and in this paper, is the *holdout estimate*. Assume we have a training sequence \mathbf{z} of length n obtained from some target concept $c \in \mathcal{H}$ by drawing inputs independently at random according to μ . We wish to make an estimate of the value of $\text{er}_\mu(L(\mathbf{z}))$ for our learner L , without making use of any further examples. Informally, the holdout estimate works as follows. We split \mathbf{z} into two sections such that $\mathbf{z} = \mathbf{z}_1\mathbf{z}_2$, use L in conjunction with \mathbf{z}_1 to obtain a hypothesis $L(\mathbf{z}_1)$, calculate the empirical error of $L(\mathbf{z}_1)$ using \mathbf{z}_2 , and use the value obtained as the required estimate. Formally, let $\mathbf{x} \in X^n$ be the sequence used to obtain \mathbf{z} . Choose an appropriate value m and divide \mathbf{x} and \mathbf{z} into m folds. The value $\text{est}_H^m(L, \mathbf{z})$ of the holdout estimate is

$$\text{est}_H^m(L, \mathbf{z}) = \text{er}_{\mathbf{x}_m}(L(\mathbf{z} \setminus \mathbf{z}_m)).$$

Note that it is not usually the case that in calculating the holdout estimate in practice we impose the restriction that m divides n . This restriction is imposed here to aid the theoretical analysis.

3.4 The cross-validation estimate

In order to calculate the holdout estimate we split the training sample into folds and use the final fold in conjunction with a hypothesis obtained using the first $(m - 1)$ folds to estimate $\text{er}_\mu(L(\mathbf{z}))$. The m -fold cross-validation estimate takes this idea a step further in that, rather than using L once to obtain a single hypothesis, it uses L once for each fold to obtain m —possibly different—hypotheses, each of which is used to obtain an estimate. These individual estimates are then averaged to obtain a final estimate. Formally, the value $\text{est}_{CV}^m(L, \mathbf{z})$ of the m -fold cross-validation estimate is given by

$$\text{est}_{CV}^m(L, \mathbf{z}) = \text{er}_{\mathbf{x}}(L(\mathbf{z} \setminus \mathbf{z}_1), \dots, L(\mathbf{z} \setminus \mathbf{z}_m)) = \frac{1}{m} \sum_{i=1}^m \text{er}_{\mathbf{x}_i}(L(\mathbf{z} \setminus \mathbf{z}_i)).$$

3.5 Previous bounds

Holden [15] obtains results in probability theory. He then applies these to obtain, for both the holdout estimate and the cross-validation estimate, upper bounds on the probability of obtaining a training sequence \mathbf{z} for which the estimate differs from the true error $\text{er}_\mu(L(\mathbf{z}))$ by more than a specified constant $0 < \epsilon \leq 1$, where L is *any*

consistent learning algorithm. Specifically, he obtains functions $f_1(n, m, \epsilon, \text{VCdim}(\mathcal{H}))$ and $f_2(n, m, \epsilon, \text{VCdim}(\mathcal{H}))$ such that under quite general conditions,

$$\mu^n \{ \mathbf{x} \in X^n : |\text{est}_H^m(L, \mathbf{z}) - \text{er}_\mu(L(\mathbf{z}))| > \epsilon \} \leq f_1(n, m, \epsilon, \text{VCdim}(\mathcal{H}))$$

and

$$\mu^n \{ \mathbf{x} \in X^n : |\text{est}_{CV}^m(L, \mathbf{z}) - \text{er}_\mu(L(\mathbf{z}))| > \epsilon \} \leq f_2(n, m, \epsilon, \text{VCdim}(\mathcal{H}))$$

where \mathbf{z} is obtained from \mathbf{x} using the target concept c . In each case, if $\text{VCdim}(\mathcal{H})$ is finite, the upper bound can be made arbitrarily small by choosing a large enough n . Also, the upper bounds given are independent of the target concept c and the distribution μ .

The main results in [15] may be stated as follows.

Theorem 3 (Holden [15]) *For any $n \geq m \geq 3$ where m divides n and $n \geq 2m/\epsilon^2$, any $c \in \mathcal{H}$, any distribution μ , and any ϵ such that $0 < \epsilon \leq 1$, let*

$$\mathcal{P}_H = \mu^n \{ \mathbf{x} \in X^n : \exists h_m, h \in \mathcal{H} \text{ with } \text{er}_{\mathbf{z} \setminus \mathbf{z}_m}(h_m) = \text{er}_{\mathbf{z}}(h) = 0, \text{er}_P(h) > \text{er}_{\mathbf{x}_m}(h_m) + \epsilon \}$$

and

$$\mathcal{P}_{CV} = \mu^n \{ \mathbf{x} \in X^n : \exists h_1, \dots, h_m, h \in \mathcal{H} \text{ with } \text{er}_{\mathbf{z} \setminus \mathbf{z}_i}(h_i) = 0 \ (1 \leq i \leq m), \text{er}_{\mathbf{z}}(h) = 0 \text{ and}$$

$$\text{er}_\mu(h) > \frac{1}{m} \sum_{i=1}^m \text{er}_{\mathbf{z}_i}(L(\mathbf{z} \setminus \mathbf{z}_i)) + \epsilon \}.$$

Then

$$\mathcal{P}_H < 2\Pi_{\mathcal{H}}^2 \left(n \left(1 + \frac{1}{m} \right) \right) 2^{-n\epsilon/2m},$$

and $\mathcal{P}_{CV} \leq m\mathcal{P}_H$.

These probability bounds imply the following results.

Theorem 4 (Holden [15]) *For any $n \geq m \geq 3$ where m divides n and $n \geq 2m/\epsilon^2$, any $c \in \mathcal{H}$, any distribution μ , and any ϵ such that $0 < \epsilon \leq 1$, we have*

$$\mu^n \{ \mathbf{x} \in X^n : |\text{est}_H^m(L, \mathbf{z}) - \text{er}_\mu(L(\mathbf{z}))| > \epsilon \} \leq 2\Pi_{\mathcal{H}}^2 \left(n \left(1 + \frac{1}{m} \right) \right) 2^{-n\epsilon/2m}$$

where $\Pi_{\mathcal{H}}^2(\cdot)$ denotes $[\Pi_{\mathcal{H}}(\cdot)]^2$, and

$$\mu^n \{ \mathbf{x} \in X^n : |\text{est}_{CV}^m(L, \mathbf{z}) - \text{er}_\mu(L(\mathbf{z}))| > \epsilon \} < 2m \Pi_{\mathcal{H}}^2 \left(n \left(1 + \frac{1}{m} \right) \right) 2^{-n\epsilon/2m}.$$

Using Sauer’s Lemma, these bounds can be translated into bounds of the form f_1, f_2 described above; see Holden [15] for details.

From this we have the following estimation error bounds.

Theorem 5 *With the notation as above, if $\delta \in (0, 1)$ is a constant, and $d = \text{VCdim}(\mathcal{H})$ is finite and $d \geq 1$, then for any consistent learning algorithm L , any $c \in \mathcal{H}$ and any μ , we have:*

1. *with probability at least $1 - \delta$,*

$$\text{er}_\mu(L(\mathbf{z})) < \text{est}_H^m(L, \mathbf{z}) + \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \left(\log \left(\frac{2}{\delta} \right) + 2 \log (\Pi_{\mathcal{H}}(n + n/m)) \right) \right\},$$

2. *with probability at least $1 - \delta$,*

$$\text{er}_\mu(L(\mathbf{z})) < \text{est}_{CV}^m(L, \mathbf{z}) + \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \left(\log \left(\frac{2m}{\delta} \right) + 2 \log (\Pi_{\mathcal{H}}(n + n/m)) \right) \right\}.$$

Sauer’s Lemma shows that if the VC-dimension of \mathcal{H} is finite, then the estimation error can be made arbitrarily small by choosing n large.

4 Using real-valued function classes

4.1 Measuring error with a margin

We now consider how one might use a class of real-valued functions for binary classification, and how we might apply a type of cross-validation to estimate the classification error. Our aim is to obtain theorems like those described above. The techniques used are similar to those used by Holden [15] and Anthony and Bartlett [5, 4].

As before, X will denote the set of all possible *examples* to be classified. (For instance, X would be \mathbf{R}^k in the case of a neural network having k real inputs.) Z will denote $X \times \{-1, 1\}$, the set of all possible *labelled examples*. We work in the general framework in which there is some unknown probability distribution P on Z , representing the concept to be learned. (As noted earlier, such a set-up includes as a special case the realisable framework in which there is a $\{-1, 1\}$ -valued target *concept* $c \in \mathcal{H}$ on X

and a probability distribution μ on X .) We consider the case where a fixed class of functions \mathcal{H} mapping into the real interval $[-1, 1]$ is used for binary classification of the examples: given $f \in \mathcal{H}$ and $x \in X$, the resulting binary classification of x by f is simply $\text{sgn}(f(x))$, where $\text{sgn}(a) = 1$ if $a \geq 0$ and $\text{sgn}(a) = -1$ if $a < 0$, and the *error* of f (with respect to P) is

$$\text{er}_P(f) = P(\{(x, y) \in X \times \{-1, 1\} : \text{sgn}(f(x)) \neq y\}).$$

The use of real-valued functions for binary classification has been considered in [22, 8, 9, 4] within versions of the PAC model of learning and versions of ‘agnostic’ PAC learning (see Kearns et al. [17] and Haussler [14]), and it has been shown that there are advantages in considering the values of the real function during training, rather than merely its sign. In particular, as suggested in [25, 13], classification with a large ‘margin’ (the distance between the classification boundary and classified examples) appears to lead in many cases to better estimates of error. Motivated by this, we study here a model of cross-validation applicable to binary classification by real-valued functions, in which one takes into account the size of the margin.

For $\gamma \in (0, 1)$, the *error of f at margin γ* (with respect to P) is

$$\text{er}_P^\gamma(f) = P(\{(x, y) \in Z : yf(x) < \gamma\}),$$

the probability that if $y = -1$ then $f(x) > -\gamma$, and if $y = 1$ then $f(x) < \gamma$. Notice that if $\text{sgn}(f(x)) \neq y$ then (remembering that $y \in \{-1, 1\}$), we certainly have $yf(x) < 0 < \gamma$, so for any $\gamma \in (0, 1)$, and any $f \in H$, $\text{er}_P^\gamma(f) \geq \text{er}_P(f)$. It will also be useful to observe that if $\gamma < \gamma'$ then $\text{er}_P^\gamma(f) \leq \text{er}_P^{\gamma'}(f)$.

The *empirical error* $\text{er}_z(f)$ of f on $\mathbf{z} \in Z^n$ is defined to be the proportion of examples in \mathbf{z} sign-misclassified by f ; that is, if $z_i = (x_i, y_i) \in X \times \{-1, 1\}$, then

$$\text{er}_z(f) = \frac{1}{n} |\{i : \text{sgn}(f(x_i)) \neq y_i\}|.$$

4.2 Measuring cross-validation error

A straightforward way in which to extend the preceding theory would be to use this measure of empirical error in place of the previous definition of empirical error. It is clear that results analogous to Theorems 4 and 5 would hold, with the VC-dimension of \mathcal{H} replaced by the VC-dimension of $\{\text{sgn}(h) : h \in \mathcal{H}\}$, where $\text{sgn}(h)(x) = \text{sgn}(h(x))$. However, we take a different approach in which we take into consideration not simply whether h gives the correct sign on a particular example, but the actual real value it gives.

Suppose that $\gamma \in (0, 1)$ is a constant. We start by consider cross-validation for algorithms which correctly classify with a margin of at least γ . Thus, such algorithms ‘interpolate’ to accuracy $1 - \gamma$. (Later, we consider more general algorithms.) Explicitly, suppose that, given a sample $\mathbf{z} \in Z^n = (X \times \{-1, 1\})^n$, our learning algorithm L produces a function $f \in H$ such that for each (x_i, y_i) in \mathbf{z} , $y_i f(x_i) \geq \gamma$. (Thus, not only does $f(x_i)$ have the correct sign on each of the examples, but the correct classification is achieved with a margin of at least γ . Note that this is, in a sense, a generalisation of consistency as defined for the standard cross-validation model, but is not as demanding as requiring that $f(x_i) = y_i$, which would be inappropriately strong.) We call such an algorithm a γ -margin algorithm. Consider how we might attempt to generalise the definition of the holdout estimate to apply to such algorithms. For a labelled sample \mathbf{z} , let \mathbf{z}_m and $\mathbf{z} \setminus \mathbf{z}_m$ be defined in the obvious way (in analogy with the definitions for the standard holdout estimate). One approach to cross-validation would be to use the empirical error $\text{er}_{\mathbf{z}_m}(L(\mathbf{z} \setminus \mathbf{z}_m))$ as an error estimate. However, we shall find some advantage in using a ‘margin-based’ estimate of the error on \mathbf{z}_m . For a constant $\gamma \in (0, 1)$, $f \in \mathcal{H}$ and $\mathbf{z} \in Z^n$, the *empirical error at margin γ* is

$$\text{er}_{\mathbf{z}}^{\gamma}(f) = \frac{1}{n} |\{i : y_i h(x_i) < \gamma\}|.$$

We shall see that in many cases, with high probability, the empirical error of $L(\mathbf{z} \setminus \mathbf{z}_m)$ at margin γ , on \mathbf{z}_m is a good estimate of the true error $\text{er}_P(L(\mathbf{z}))$ of L on the whole sample \mathbf{z} . That is, we shall obtain results concerning the estimator

$$\text{est}_{H, \gamma}^m(L, \mathbf{z}) = \text{er}_{\mathbf{z}_m}^{\gamma}(L(\mathbf{z} \setminus \mathbf{z}_m)).$$

As in [15], we shall relate the cross-validation estimate to the holdout estimate and thereby obtain a similar result for the corresponding cross-validation estimate

$$\text{est}_{CV, \gamma}^m(L, \mathbf{z}) = \frac{1}{m} \sum_{i=1}^m \text{er}_{\mathbf{z}_i}^{\gamma}(L(\mathbf{z} \setminus \mathbf{z}_i)).$$

5 Results for real-valued cross-validation models

5.1 Covering numbers

To state our main result, we require the important notion of *covering number*. Given $W \subseteq \mathbf{R}^k$ and a positive real number ϵ , we say that $C \subseteq \mathbf{R}^k$ is a d_{∞} ϵ -cover for W if $C \subseteq W$ and for every $w \in W$ there is a $v \in C$ such that

$$\max \{|w_i - v_i| : i = 1, \dots, k\} < \epsilon.$$

We say that W is *totally bounded* if for each $\epsilon > 0$ there is a *finite* ϵ -cover for W . In this case, we define the d_∞ ϵ -covering number of W , $\mathcal{N}(W, \epsilon, d_\infty)$, to be the minimum cardinality of a d_∞ ϵ -cover for W .

Suppose that F is a class of functions from X to \mathbf{R} . Given $\mathbf{x} = (x_1, x_2, \dots, x_k) \in X^k$, we let $F|_{\mathbf{x}}$ be the subset of \mathbf{R}^k given by

$$F|_{\mathbf{x}} = \{(f(x_1), f(x_2), \dots, f(x_k)) : f \in F\}.$$

For a positive number ϵ , we define the (*uniform*) covering number $\mathcal{N}_\infty(F, \epsilon, k)$ to be the maximum, over all $\mathbf{x} \in X^k$, of the covering number $\mathcal{N}(F|_{\mathbf{x}}, \epsilon, d_\infty)$ (and we take it to be infinite if there is no bound on these covering numbers); that is,

$$\mathcal{N}_\infty(F, \epsilon, k) = \max \left\{ \mathcal{N}(F|_{\mathbf{x}}, \epsilon, d_\infty) : \mathbf{x} \in X^k \right\}.$$

The covering number is a generalisation of the growth function. To see this, suppose that functions in H map into $\{-1, 1\}$. Then, for all $\mathbf{x} \in X^k$, $H|_{\mathbf{x}}$ is finite and, for all $\epsilon < 1$, $\mathcal{N}(H|_{\mathbf{x}}, \epsilon, d_\infty) = |H|_{\mathbf{x}}|$, so $\mathcal{N}_\infty(H, \epsilon, k) = \Pi_H(k)$.

5.2 Main probability theorems

Our main theorem for the form of holdout estimate defined above is as follows. Note that this result is concerned with showing that $\text{est}_{H, \gamma}^m(L, \mathbf{z})$ does not *underestimate* by too much the true error $\text{er}_P(L(\mathbf{z}))$ of $L(\mathbf{z})$. It does not concern the probability that $\text{est}_{H, \gamma}^m(L, \mathbf{z})$ overestimates the true error. That is, we bound the probability that $\text{er}_P(L(\mathbf{z})) > \text{est}_{H, \gamma}^m(L, \mathbf{z}) + \epsilon$, rather than the probability that $|\text{er}_P(L(\mathbf{z})) - \text{est}_{H, \gamma}^m(L, \mathbf{z})| > \epsilon$. (If, in fact, $\text{er}_P(L(\mathbf{z})) < \text{est}_{H, \gamma}^m(L, \mathbf{z}) - \epsilon$ then using the error estimate does not, in any case, mislead us into thinking that a bad hypothesis is in fact a good one.)

Theorem 6 *Let $\gamma \in (0, 1)$. For $n \geq m \geq 3$, where m divides n and $n \geq 2m/\epsilon^2$, for any distribution P on $Z = X \times \{-1, 1\}$, and any $\epsilon \in (0, 1]$, let*

$$\mathcal{P}_H = P^n \left(\left\{ \mathbf{z} \in Z^n : \exists h_m, h \in \mathcal{H} \text{ with } \text{er}_{\mathbf{z} \setminus \mathbf{z}_m}^\gamma(h_m) = 0, \text{er}_{\mathbf{z}}^\gamma(h) = 0 \text{ and } \text{er}_P(h) \geq \text{er}_{\mathbf{x}_m}^\gamma(h_m) + \epsilon \right\} \right).$$

Then

$$\mathcal{P}_H < 2 (\mathcal{N}_\infty(\mathcal{H}, \gamma/2, n + n/m))^2 2^{-n\epsilon/2m}.$$

Equivalently, for $\delta, \gamma \in (0, 1)$ and for n, m as above, with probability at least $1 - \delta$, if $h_m, h \in \mathcal{H}$ and $\text{er}_{\mathbf{z} \setminus \mathbf{z}_m}^\gamma(h_m) = \text{er}_{\mathbf{z}}^\gamma(h) = 0$ then

$$\text{er}_P(h) < \text{er}_{\mathbf{x}_m}^\gamma(h_m) + \epsilon(n, \delta, \gamma),$$

where

$$\epsilon(n, m, \delta, \gamma) = \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{2}{\delta} \right) + \frac{4m}{n} \log (\mathcal{N}_\infty (\mathcal{H}, \gamma/2, n + n/m)) \right\}.$$

In particular, therefore, if L is a γ -margin algorithm, we have that, with probability at least $1 - \delta$,

$$\text{er}_P(L(\mathbf{z})) < \text{est}_{H,\gamma}^m(L, \mathbf{z}) + \epsilon(n, m, \delta, \gamma).$$

For the modified cross-validation estimate, we have the following result.

Theorem 7 For $n \geq m \geq 3$, where m divides n and $n \geq 2m/\epsilon^2$, and with the notation as above, if

$$\mathcal{P}_{\text{CV}} = P^n \left(\left\{ \mathbf{z} \in Z^n : \exists h_1, \dots, h_m, h \in \mathcal{H} \text{ with } \text{er}_{\mathbf{z} \setminus \mathbf{z}_i}^\gamma(h_i) = 0 \ (1 \leq i \leq m), \text{er}_{\mathbf{z}}^\gamma(h) = 0 \text{ and } \text{er}_P(h) > \frac{1}{m} \sum_{i=1}^m \text{er}_{\mathbf{z}_i}^\gamma(h_i) + \epsilon \right\} \right),$$

then

$$\mathcal{P}_{\text{CV}} < 2m (\mathcal{N}_\infty (\mathcal{H}, \gamma/2, n + n/m))^2 2^{-n\epsilon/2m}.$$

Equivalently, for $\delta, \gamma \in (0, 1)$ and for n, m as above, with probability at least $1 - \delta$, if $h_1, \dots, h_m, h \in \mathcal{H}$ and $\text{er}_{\mathbf{z} \setminus \mathbf{z}_i}^\gamma(h_i) = 0$ for $1 \leq i \leq m$, and $\text{er}_{\mathbf{z}}^\gamma(h) = 0$ then

$$\text{er}_P(h) < \frac{1}{m} \sum_{i=1}^m \text{er}_{\mathbf{z}_i}^\gamma(h_i) + \epsilon(n, m, \delta, \gamma),$$

where

$$\epsilon(n, m, \delta, \gamma) = \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{2m}{\delta} \right) + \frac{4m}{n} \log (\mathcal{N}_\infty (\mathcal{H}, \gamma/2, n + n/m)) \right\}.$$

Hence, for any γ -margin algorithm L , we have that, with probability at least $1 - \delta$,

$$\text{er}_P(L(\mathbf{z})) < \text{est}_{\text{CV},\gamma}^m(L, \mathbf{z}) + \epsilon'(n, m, \delta, \gamma).$$

It is perhaps unnatural to specify in advance the parameter γ . Indeed, it would seem more appropriate to compute γ on the basis of the training sample, and this will be the basis of the algorithm to be discussed later, in Section 6 (in which we define ‘sample-based’ modified holdout and cross-validation estimators). The following results will be useful.

Theorem 8 For $\delta \in (0, 1)$ and for $n \geq m \geq 3$, where m divides n , with probability at least $1 - \delta$, for all $\gamma \in (0, 1)$, if $h_m, h \in \mathcal{H}$ and $\text{er}_{\mathbf{z} \setminus \mathbf{z}_m}^\gamma(h_m) = \text{er}_{\mathbf{z}}^\gamma(h) = 0$ then

$$\text{er}_P(h) < \text{er}_{\mathbf{x}_m}^\gamma(h_m) + \varepsilon(n, m, \delta, \gamma),$$

where

$$\varepsilon(n, m, \delta, \gamma) = \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{4}{\delta\gamma} \right) + \frac{4m}{n} \log (\mathcal{N}_\infty (\mathcal{H}, \gamma/4, n + n/m)) \right\}.$$

Theorem 9 For $\delta \in (0, 1)$ and for $n \geq m \geq 3$, where m divides n , with probability at least $1 - \delta$, for all $\gamma \in (0, 1)$, if $h_1, \dots, h_m, h \in \mathcal{H}$ and $\text{er}_{\mathbf{z} \setminus \mathbf{z}_i}^\gamma(h_i) = 0$ for $1 \leq i \leq m$, and $\text{er}_{\mathbf{z}}^\gamma(h) = 0$, then

$$\text{er}_P(h) < \frac{1}{m} \sum_{i=1}^m \text{er}_{\mathbf{z}_i}^\gamma(h_i) + \varepsilon(n, m, \delta, \gamma),$$

where

$$\varepsilon(n, m, \delta, \gamma) = \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{4m}{\delta\gamma} \right) + \frac{4m}{n} \log (\mathcal{N}_\infty (\mathcal{H}, \gamma/4, n + n/m)) \right\}.$$

5.3 Discussion

Comparison with the binary cross-validation bounds

If we ignore the margin of classification and simply use the sign $\text{sgn}(h)$ of $f : X \rightarrow [-1, 1]$, then Theorem 4 will provide bounds on the hold-out and cross-validation estimates, with \mathcal{H} relaxed by $\text{sgn}(\mathcal{H}) = \{\text{sgn}(h) : h \in \mathcal{H}\}$. These bounds will involve the growth function of the class $\text{sgn}(\mathcal{H})$. If $\text{sgn}(\mathcal{H})$ has finite VC-dimension then, by Sauer's Lemma, this growth function will be subexponential and the estimation error can be made arbitrarily small by choosing n large enough. On the other hand, taking account of the margin and using the results just given, then the relevant bounds involve the covering numbers. In many cases, these are subexponential even when $\text{sgn}(\mathcal{H})$ has infinite dimension. (This will happen if $\text{sgn}(\mathcal{H})$ has infinite VC-dimension but \mathcal{H} has finite fat-shattering dimension—see later.) Thus the new results provide useful bounds on the error in situations where the previous bounds tell us nothing at all.

Sanity-checking and existential anxiety

Comparing our new bounds to ones which apply to an appropriate resubstitution error, however, makes it clear that our bounds are ‘sanity-check’ bounds, to use the termi-

nology of Kearns and Ron [18]. Explicitly, it has been shown in [22, 8, 9] that the following holds for $\gamma, \delta \in (0, 1)$: with probability at least $1 - \delta$, if $\text{er}_{\mathbf{z}}^{\gamma}(h) = 0$ then

$$\text{er}_P(h) < \frac{4}{n} \ln \left(\frac{4}{\delta} \right) + \frac{4}{n} \ln (\mathcal{N}_{\infty}(\mathcal{H}, \gamma/2, 2n)).$$

Thus, the γ -resubstitution estimate $\text{er}_{\mathbf{z}}^{\gamma}(h)$ has a smaller estimation error than we have been able to obtain for the cross-validation estimator. The same is true of the binary-valued case (see Holden [15]). The looseness in our bound is a result of the way we derive it from a probability bound. Explicitly, we bound the probability \mathcal{P}_{CV} of the set of $\mathbf{z} \in Z^n$ such that $\exists h_1, \dots, h_m, h \in \mathcal{H}$ with $\text{er}_{\mathbf{z} \setminus \mathbf{z}_i}^{\gamma}(h_i) = 0$ for $1 \leq i \leq m$, $\text{er}_{\mathbf{z}}^{\gamma}(h) = 0$, and

$$\text{er}_P(h) > \frac{1}{m} \sum_{i=1}^m \text{er}_{\mathbf{z}_i}^{\gamma}(h_i) + \epsilon.$$

To obtain a bound for cross-validation, the existential quantifiers are not needed: that is, instead of asking whether there *exists* h_i and h with the required property, we more specifically question whether the particular functions $L(\mathbf{z} \setminus \mathbf{z}_i)$ and $L(\mathbf{z})$ have the property. That is to say, the functions h_1, h_2, \dots, h_m, h are not independently chosen but arise from L . Now, if $g \in \mathcal{H}$ is such that $\text{er}_{\mathbf{z}}^{\gamma}(g) = 0$ and $\text{er}_P(g) > \epsilon$, then, taking each h_i , and h , to equal g , we see that \mathbf{z} belongs to the set just described. Thus,

$$\mathcal{P}_{\text{CV}} \geq P^n \{ \mathbf{z} : \exists g \in \mathcal{H}, \text{er}_{\mathbf{z}}^{\gamma}(g) = 0, \text{er}_P(g) > \epsilon \}.$$

The estimation error bounds for the γ -resubstitution error follow by bounding the right-hand probability from above, so it follows that, using the general method of this paper, it is impossible to obtain bounds which show that cross-validation is better than resubstitution. (The same is true of the holdout estimate.) It might therefore be argued that the approach we take is over-general and that a different approach would yield better than mere ‘sanity-check’ bounds. This is an open problem. We suggest, however, that our probability bounds are of interest in their own right. In any case, our results are of use in that they provide reasonable bounds for cross-validation in situations where no previous bounds were known.

6 New estimators

6.1 Modified estimators

Theorem 8 and Theorem 9 allow one to compute a suitable γ on the basis of the performance of the learning algorithm in the training sample. Consider for the moment the hold-out estimate. Suppose L is a learning algorithm and that a sample \mathbf{z} is given.

Let $\gamma(\mathbf{z})$ be the largest γ for which

$$\text{er}_{\mathbf{z} \setminus \mathbf{z}_m}(L(\mathbf{z} \setminus \mathbf{z}_m)) = \text{er}_{\mathbf{z}}(L(\mathbf{z})) = 0.$$

Then, using this value of γ , and applying Theorem 8, we have a high-probability bound on the true error of $L(\mathbf{z})$: explicitly, $\text{er}_P(L(\mathbf{z}))$ is no more than $\varepsilon(n, m, \delta, \gamma)$ plus $\text{er}_{\mathbf{z}_m}^\gamma(L(\mathbf{z} \setminus \mathbf{z}_m))$, the error of $L(\mathbf{z} \setminus \mathbf{z}_m)$ at margin γ on \mathbf{z}_m .

For $\mathbf{w} = ((x_1, y_1), \dots, (x_k, y_k)) \in Z^k$ and $h \in H$, let

$$\text{mar}(f, \mathbf{w}) = \min\{f(x_i)y_i : 1 \leq i \leq k\}.$$

We formalise the above as follows.

Definition 10 (Modified holdout estimate) For $n \geq m \geq 1$, for a learning algorithm L , and for $\mathbf{z} \in Z^n$, let

$$\text{Hmar}(L, \mathbf{z}) = \min\{\text{mar}(L(\mathbf{z} \setminus \mathbf{z}_m), \mathbf{z} \setminus \mathbf{z}_m), \text{mar}(L(\mathbf{z}), \mathbf{z})\}.$$

Then the modified holdout estimate $\text{Mest}_H^m(L, \mathbf{z})$ is defined to be 1 if $\text{Hmar}(L, \mathbf{z}) \leq 0$, and otherwise is given by

$$\text{Mest}_H^m(L, \mathbf{z}) = \text{est}_{H, \text{Hmar}(L, \mathbf{z})}^m(L, \mathbf{z}) = \text{er}_{\mathbf{z}_m}^{\text{Hmar}(L, \mathbf{z})}(L(\mathbf{z} \setminus \mathbf{z}_m)).$$

In a similar way, we define a modified cross-validation estimator.

Definition 11 (Modified cross-validation estimate) For $n \geq m \geq 1$, for a learning algorithm L , and for $\mathbf{z} \in Z^n$, let

$$\text{CVmar}(L, \mathbf{z}) = \min\left\{\min_{1 \leq i \leq m} \text{mar}(L(\mathbf{z} \setminus \mathbf{z}_i), \mathbf{z} \setminus \mathbf{z}_i), \text{mar}(L(\mathbf{z}), \mathbf{z})\right\}.$$

Then the modified cross-validation estimate $\text{Mest}_{\text{CV}}^m(L, \mathbf{z})$ is defined to be 1 if $\text{CVmar}(L, \mathbf{z}) \leq 0$, and otherwise is given by

$$\text{Mest}_{\text{CV}}^m(L, \mathbf{z}) = \text{est}_{H, \text{CVmar}(L, \mathbf{z})}^m(L, \mathbf{z}) = \frac{1}{m} \sum_{i=0}^m \text{er}_{\mathbf{z}_m}^{\text{CVmar}(L, \mathbf{z})}(L(\mathbf{z} \setminus \mathbf{z}_i)).$$

6.2 Performance bounds

Using Theorem 8 and Theorem 9, we obtain the following bounds on the error of the modified estimators.

Theorem 12 For $\delta \in (0, 1)$ and for $n \geq m \geq 3$, where m divides n , with probability at least $1 - \delta$,

1. if $\text{Hmar}(L, \mathbf{z}) > 0$, then

$$\text{er}_P(L(\mathbf{z})) < \text{Mest}_H^m(L, \mathbf{z}) + \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{4}{\delta\gamma} \right) + \frac{4m}{n} \log (\mathcal{N}_\infty (\mathcal{H}, \gamma/4, n + n/m)) \right\},$$

where $\gamma = \text{Hmar}(L, \mathbf{z})$.

2. if $\text{CVmar}(L, \mathbf{z}) > 0$, then

$$\text{er}_P(L(\mathbf{z})) < \text{Mest}_{\text{CV}}^m(L, \mathbf{z}) + \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{4m}{\delta\gamma} \right) + \frac{4m}{n} \log (\mathcal{N}_\infty (\mathcal{H}, \gamma/4, n + n/m)) \right\},$$

where $\gamma = \text{CVmar}(L, \mathbf{z})$.

6.3 Application to neural networks

We now use bounds of Bartlett [8, 9, 4] on the covering numbers of function classes computable by certain feedforward sigmoid neural networks. We deal here with networks having as activation function $\sigma(y) = \tanh(y/2) = 1 - 1/(1 + e^{-y})$ and, for convenience, we assume that all thresholds are set at 0. (There are similar results for other types of sigmoid function, and for nets with thresholds, but to illustrate the results as concisely as possible, we shall make these assumptions.) We define

$$F_0 = \{x \mapsto x_i : x = (x_1, \dots, x_n) \in [-1, 1]^r, i \in \{1, \dots, r\}\},$$

and, for $i \geq 1$, we let

$$F_i = \left\{ \sigma \left(\sum_{j=1}^N w_j f_j \right) : N \in \mathbf{N}, f_j \in \bigcup_{k=0}^{i-1} F_k, \sum_{j=1}^N |w_j| \leq V \right\}.$$

Thus, $\mathcal{H} = F_\ell$ is the class of functions that can be computed by an ℓ -layer feed-forward network in which each unit has the sum of the magnitudes of its weights bounded by V .

Bartlett has shown ([4], following [8, 9]) that

$$\log_2 \mathcal{N}_2 (F_\ell, \epsilon, m) \leq \frac{1}{6} \left(\frac{3}{\epsilon} \right)^{2\ell} (2V)^{\ell(\ell+1)} \log(2r + 2),$$

provided $V \geq 1/2$, and $\epsilon \leq V$.

We therefore have the following result.

Theorem 13 Let $\delta \in (0, 1)$ and suppose $n \geq m \geq 3$, where m divides n . Let \mathcal{H} be the class of functions computable by the type of neural network just described (with weight-bound V and r inputs), and assume that $V \geq 1/2$. Then, with probability at least $1 - \delta$:

1. if $\text{Hmar}(L, \mathbf{z}) > 0$, then

$$\text{er}_P(L(\mathbf{z})) < \text{Mest}_H^m(L, \mathbf{z}) + \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{4}{\delta\gamma} \right) + \frac{2m}{3n} \left(\frac{12}{\gamma} \right) (2V)^{\ell(\ell+1)} \log(2r+2) \right\},$$

where $\gamma = \min\{\text{Hmar}(L, \mathbf{z}), 4V\}$;

2. if $\text{CVmar}(L, \mathbf{z}) > 0$, then

$$\text{er}_P(L(\mathbf{z})) < \text{Mest}_{CV}^m(L, \mathbf{z}) + \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{4m}{\delta\gamma} \right) + \frac{2m}{3n} \left(\frac{12}{\gamma} \right) (2V)^{\ell(\ell+1)} \log(2r+2) \right\},$$

where $\gamma = \min\{\text{CVmar}(L, \mathbf{z}), 4V\}$.

7 General ‘dimension-based’ bounds

In this section, we return to the general theory. The general bounds obtained earlier concern the covering numbers of the function class \mathcal{H} . It is often useful to have bounds in terms of certain dimensions associated with the function class, analogous to the VC-dimension for binary-valued classes.

A number of ways of measuring the ‘expressive power’ of a class \mathcal{H} of functions have been found to be useful. This power can be quantified by associating a ‘dimension’ to the class. Sometimes this is simply one number depending on \mathcal{H} and sometimes—in what is known as a *scale-sensitive dimension*—it is a function depending on \mathcal{H} .

An important example of the first type of dimension is the *pseudo-dimension* [14, 20]. We say that a finite subset $S = \{x_1, x_2, \dots, x_d\}$ of X is *shattered* if there is $\mathbf{r} = (r_1, r_2, \dots, r_d) \in \mathbf{R}^d$ such that for every $b = (b_1, b_2, \dots, b_d) \in \{0, 1\}^d$, there is a function $h_b \in \mathcal{H}$ with $h_b(x_i) > r_i$ if $b_i = 1$ and $h_b(x_i) < r_i$ if $b_i = 0$. The *pseudo-dimension* of \mathcal{H} , denoted $\text{Pdim}(\mathcal{H})$, is the largest cardinality of a shattered set, or infinity if there is no bound on the cardinalities of the shattered sets. The pseudo-dimension is a well-understood and useful measure of expressive power. One attractive feature of this dimension is that if the set of functions is a vector space then its pseudo-dimension coincides with its linear dimension (see for example [14]).

The most important scale-sensitive dimension which has been used to date in the development of the theory of learning real-valued functions is the *fat-shattering dimension*. This is a scale-sensitive version of the pseudo-dimension and was introduced by Kearns and Schapire [19]. Suppose that \mathcal{H} is a set of functions from X to $[-1, 1]$ and that $\gamma \in (0, 1)$. We say that a finite subset $S = \{x_1, x_2, \dots, x_d\}$ of X is γ -shattered if there is $\mathbf{r} = (r_1, r_2, \dots, r_d) \in \mathbf{R}^d$ such that for every $b = (b_1, b_2, \dots, b_d) \in \{0, 1\}^d$, there is a function $h_b \in \mathcal{H}$ with $h_b(x_i) \geq r_i + \gamma$ if $b_i = 1$ and $h_b(x_i) \leq r_i - \gamma$ if $b_i = 0$. Thus, S is γ -shattered if it is shattered with a ‘width of shattering’ of at least γ . We define the *fat-shattering dimension*, $\text{fat}_{\mathcal{H}} : \mathbf{R}^+ \rightarrow \mathbf{N}_0 \cup \{\infty\}$, as

$$\text{fat}_{\mathcal{H}}(\gamma) = \max \{|S| : S \subseteq X \text{ is } \gamma\text{-shattered by } \mathcal{H}\},$$

or $\text{fat}_{\mathcal{H}}(\gamma) = \infty$ if the maximum does not exist. (Here, \mathbf{N}_0 denotes the set of nonnegative integers.) It is easy to see that $\text{Pdim}(\mathcal{H}) = \lim_{\gamma \rightarrow 0} \text{fat}_{\mathcal{H}}(\gamma)$. It should be noted, however, that it is possible for the pseudo-dimension to be infinite, even when $\text{fat}_{\mathcal{H}}(\gamma)$ is finite for all γ . We shall say that \mathcal{H} has *finite fat-shattering dimension* whenever it is the case that for all $\gamma \in (0, 1)$, $\text{fat}_{\mathcal{H}}(\gamma)$ is finite.

The following result is due to Alon *et al.* [6]. It bounds the d_{∞} -covering number of \mathcal{H} in terms of the fat-shattering function of \mathcal{H} .

Theorem 14 *Suppose that \mathcal{H} is a set of real functions from a domain X to the bounded interval $[-1, 1]$ and that $\epsilon > 0$. Then*

$$\mathcal{N}_{\infty}(\mathcal{H}, \epsilon, m) < 2 \left(mb^2 \right)^{\lceil \log_2 y \rceil} < 2 \left(\frac{16m}{\epsilon^2} \right)^{d \log_2(4em/(d\epsilon))},$$

where

$$b = \left\lfloor \frac{4}{\epsilon} \right\rfloor,$$

and, with $d = \text{fat}_{\mathcal{H}}(\epsilon/4)$,

$$y = \sum_{i=1}^d \binom{m}{i} b^i.$$

Using Theorem 8 and Theorem 9 we obtain the following bounds in terms of the fat-shattering dimension. (The proofs are omitted.)

Theorem 15 *For $\delta \in (0, 1)$ and for $n \geq m \geq 3$, where m divides n , with probability at least $1 - \delta$,*

1. *for all $\gamma \in (0, 1)$, if $h_m, h \in \mathcal{H}$ and $\text{er}_{\mathbf{z} \setminus \mathbf{z}_m}^{\gamma}(h_m) = \text{er}_{\mathbf{z}}^{\gamma}(h) = 0$ then*

$$\text{er}_P(h) < \text{er}_{\mathbf{x}_m}^{\gamma}(h_m) + \varepsilon(n, m, \delta, \gamma),$$

where

$$\varepsilon(n, m, \delta, \gamma) = \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{4}{\delta\gamma} \right) + \frac{16m}{n} \text{fat}_H(\gamma/16) \log^2 \left(\frac{23n}{\gamma} \right) \right\}.$$

2. for all $\gamma \in (0, 1)$, if $h_1, \dots, h_m, h \in \mathcal{H}$ and $\text{er}_{\mathbf{z} \setminus \mathbf{z}_i}^\gamma(h_i) = 0$ for $1 \leq i \leq m$, and $\text{er}_{\mathbf{z}}^\gamma(h) = 0$, then

$$\text{er}_P(h) < \frac{1}{m} \sum_{i=1}^m \text{er}_{\mathbf{z}_i}^\gamma(h_i) + \varepsilon(n, m, \delta, \gamma),$$

where

$$\varepsilon(n, m, \delta, \gamma) = \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{4m}{\delta\gamma} \right) + \frac{16m}{n} d \log^2 \left(\frac{23n}{\gamma} \right) \right\}.$$

In particular, for the modified holdout and cross-validation estimators, we have the following results.

Theorem 16 For $n \geq m \geq 3$, where m divides n and $n \geq 2m/\epsilon^2$, and with the notation as above, we have that, with probability at least $1 - \delta$,

1. if $\text{Hmar}(L, \mathbf{z}) > 0$, then

$$\text{er}_P(L(\mathbf{z})) < \text{Mest}_H^m(L, \mathbf{z}) + \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{4}{\delta\gamma} \right) + \frac{16m}{n} d \log^2 \left(\frac{23n}{\gamma} \right) \right\},$$

where $\gamma = \text{Hmar}(L, \mathbf{z})$ and where $d = \text{fat}_{\mathcal{H}}(\text{Hmar}(L, \mathbf{z})/16)$.

2. if $\text{CVmar}(L, \mathbf{z}) > 0$, then

$$\text{er}_P(L(\mathbf{z})) < \text{est}_{CV, \gamma}^m(L, \mathbf{z}) + \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{4m}{\delta\gamma} \right) + \frac{16m}{n} d \log^2 \left(\frac{23n}{\gamma} \right) \right\},$$

where $\gamma = \text{CVmar}(L, \mathbf{z})$ and $d = \text{fat}_{\mathcal{H}}(\gamma/16)$.

As indicated earlier, these results apply in many cases where the binary results applied to $\text{sgn}(\mathcal{H})$ tell us nothing. Explicitly, if \mathcal{H} has finite fat-shattering dimension but $\text{VCdim}(\text{sgn}(\mathcal{H}))$ is infinite, then the binary results are useless. (There are many such classes. For example, the set of functions $f : [-1, 1] \rightarrow [-1, 1]$ of total variation at most 1 is such a class.)

8 Proofs

8.1 Relating cross-validation estimate to holdout estimate

We start by defining, for any \mathbf{z} , the following two sets of hypothesis sequences:

$$\mathcal{H}_{H,\gamma}^m(\mathbf{z}) = \left\{ \mathbf{h} = (h_m, h) \in \mathcal{H}^2 : \text{er}_{\mathbf{z} \setminus \mathbf{z}_m}^\gamma(h_m) = 0, \text{er}_{\mathbf{z}}^\gamma(h) = 0 \right\}$$

and

$$\mathcal{H}_{CV,\gamma}^m(\mathbf{z}) = \left\{ \mathbf{h} = (h_1, \dots, h_m, h) \in \mathcal{H}^{m+1} : \text{er}_{\mathbf{z} \setminus \mathbf{z}_i}^\gamma(h_i) = 0 \text{ for } 1 \leq i \leq m \text{ and } \text{er}_{\mathbf{z}}^\gamma(h) = 0. \right\}$$

We then define the sets

$$Q_{n,m}^{\epsilon,\gamma} = \left\{ \mathbf{z} \in Z^n : \exists \mathbf{h} \in \mathcal{H}_{H,\gamma}^m(\mathbf{z}) \text{ such that } \text{er}_P(h) > \text{er}_{\mathbf{z}_m}(h_m) + \epsilon \right\}$$

and

$$R_{n,m}^{\epsilon,\gamma} = \left\{ \mathbf{z} \in Z^n : \exists \mathbf{h} \in \mathcal{H}_{CV,\gamma}^m(\mathbf{z}) \text{ such that } \text{er}_P(h) > \text{er}_{\mathbf{z}}(h_1, \dots, h_m) + \epsilon \right\}.$$

As in the corresponding analysis of [15], an upper bound for $P^n(Q_{n,m}^{\epsilon,\gamma})$ immediately implies one for $P^n(R_{n,m}^\epsilon)$, as we now show.

Lemma 17 *For any P , any $n \geq 2$, any $\gamma \in (0, 1)$, any m such that $m \geq 2$ and m divides n , and any $\epsilon \in (0, 1]$,*

$$P^n(R_{n,m}^{\epsilon,\gamma}) \leq mP^n(Q_{n,m}^{\epsilon,\gamma}).$$

Proof Recalling the definition of the set $R_{n,m}^{\epsilon,\gamma}$,

$$\begin{aligned} R_{n,m}^\epsilon &= \left\{ \mathbf{z} \in Z^n : \exists \mathbf{h} \in \mathcal{H}_{CV,\gamma}^m(\mathbf{z}) \text{ such that } \text{er}_P(h) > \text{er}_{\mathbf{z}}(h_1, \dots, h_m) + \epsilon \right\} \\ &= \left\{ \mathbf{z} \in Z^n : \exists \mathbf{h} \in \mathcal{H}_{CV,\gamma}^m(\mathbf{z}) \text{ such that } \text{er}_P(h) > \frac{1}{m} \sum_{i=1}^m \text{er}_{\mathbf{z}_i}^\gamma(L(\mathbf{z} \setminus \mathbf{z}_i)) + \epsilon \right\} \\ &\subseteq \left\{ \mathbf{z} \in Z^n : \exists \mathbf{h} \in \mathcal{H}_{CV,\gamma}^m(\mathbf{z}) \text{ s.t. } \text{er}_P(h) > \text{er}_{\mathbf{z}_i}^\gamma(L(\mathbf{z} \setminus \mathbf{z}_i)) + \epsilon \text{ for some } i, 1 \leq i \leq m \right\} \\ &\subseteq \bigcup_{i=1}^m A_{n,m}^{\epsilon,\gamma}(i) \end{aligned}$$

where the sets $A_{n,m}^{\epsilon,\gamma}(i)$ are defined as

$$A_{n,m}^{\epsilon,\gamma}(i) = \left\{ \mathbf{z} \in Z^n : \exists h, h_i \in \mathcal{H} \text{ with } \text{er}_{\mathbf{z} \setminus \mathbf{z}_i}^\gamma(h_i) = 0, \text{er}_{\mathbf{z}}^\gamma(h) = 0, \text{er}_P(h) > \text{er}_{\mathbf{z}_i}(h_i) + \epsilon \right\}$$

for $1 \leq i \leq m$. Consequently,

$$P^n(R_{n,m}^{\epsilon,\gamma}) \leq P^n \left(\bigcup_{i=1}^m A_{n,m}^{\epsilon,\gamma}(i) \right) \leq \sum_{i=1}^m P^n(A_{n,m}^{\epsilon,\gamma}(i)).$$

By symmetry, we have $P^n(A_{n,m}^{\epsilon,\gamma}(i)) = P^n(Q_{n,m}^{\epsilon,\gamma})$ for $1 \leq i \leq m$, and the required results follows. \square

8.2 Symmetrization

We will present the proof of Theorem 6 using a series of lemmas. The techniques are similar to those used in [15] and [5], which in turn extend the standard techniques of symmetrization and combinatorial bounding used in [26, 10, 3]. For notational convenience we will divide any $\mathbf{z} \in Z^n$ into two parts and write $\mathbf{z} = \mathbf{z}'\mathbf{z}_m$ where $\mathbf{z}' = \mathbf{z} \setminus \mathbf{z}_m = (z_1, \dots, z_{n-n'})$ and as usual $\mathbf{z}_m = (z_{n-n'+1}, \dots, z_n)$. We will also introduce a further sequence $\mathbf{w} = (y_1, \dots, y_k) \in X^k$, where the value of k will be fixed below. Define the set

$$S_{n,m}^{\epsilon,\gamma}(k) = \left\{ \mathbf{z}'\mathbf{z}_m\mathbf{w} \in Z^{n+k} : \exists \mathbf{h} \in \mathcal{H}_{H,\gamma}^m(\mathbf{z}) \text{ such that } \text{er}_{\mathbf{w}}(h) > \text{er}_{\mathbf{z}_m}^\gamma(h_m) + \epsilon/2 \right\}.$$

Note that the error measure $\text{er}_{\mathbf{w}}(h)$ in this definition is the (standard) empirical error; that is, the proportion of examples (x_i, y_i) in \mathbf{w} such that $\text{sgn}(h(x_i)) \neq y_i$.

Lemma 18 *For any $n \geq 2$, any m such that $m \geq 2$ and m divides n , any ϵ such that $0 < \epsilon \leq 1$, any $\gamma \in (0, 1)$, any distribution P , and any $k \geq 2/\epsilon^2$,*

$$P^n(Q_{n,m}^{\epsilon,\gamma}) \leq 2P^{n+k}(S_{n,m}^{\epsilon,\gamma}(k)).$$

Proof If $h \in \mathcal{H}$ satisfies $\text{er}_P(h) > \text{er}_{\mathbf{z}_m}^\gamma(h_m) + \epsilon$ and $\text{er}_{\mathbf{w}}(h) \geq \text{er}_P(h) - \epsilon/2$, then

$$\text{er}_{\mathbf{w}}(h) > \text{er}_{\mathbf{z}_m}^\gamma(h_m) + \epsilon/2.$$

so

$$\begin{aligned} P^{n+k}(S_{n,m}^{\epsilon,\gamma}) &\geq P^{n+k} \left(\left\{ \mathbf{z}\mathbf{w} : \exists h \in \mathcal{H}_{H,\gamma}^m(\mathbf{z}), \text{er}_P(h) \geq \text{er}_{\mathbf{z}_m}^\gamma(h_m) + \epsilon \text{ and } \text{er}_{\mathbf{w}}(h) \geq \text{er}_P(h) - \epsilon/2 \right\} \right) \\ &= \int_{Z^n} I_Q(\mathbf{z}) P^k \left(\left\{ \mathbf{w} : \exists h \in \mathcal{H}, \text{er}_P(h) \geq \text{er}_{\mathbf{z}_m}^\gamma(h_m) + \epsilon \text{ and } \right. \right. \\ &\quad \left. \left. \text{er}_{\mathbf{w}}(h) \geq \text{er}_P(h) - \epsilon/2 \right\} \right) dP^n(\mathbf{z}), \end{aligned} \tag{1}$$

where I_Q denotes the indicator (or characteristic) function of the set $Q_{n,m}^{\epsilon,\gamma}$.

Now, for $\mathbf{z} \in Q_{n,m}^{\epsilon,\gamma}$ fix an $h \in \mathcal{H}$ with $\text{er}_P(h) \geq \text{er}_{\mathbf{z}_m}^\gamma(h_m) + \epsilon$. For this h , we shall show that for $k \geq 2/\epsilon^2$,

$$P^k(\{\mathbf{w} : \text{er}_{\mathbf{w}}(h) \geq \text{er}_P(h) - \epsilon/2\}) \geq 1/2. \quad (2)$$

It will then follow that, for any $\mathbf{z} \in Q = Q_{n,m}^{\epsilon,\gamma}$, we have

$$P^k\left(\left\{\mathbf{w} : \exists h \in \mathcal{H}, \text{er}_P(h) \geq \text{er}_{\mathbf{z}_m}^\gamma(h_m) + \epsilon \text{ and } \text{er}_{\mathbf{w}}(h) \geq \text{er}_P(h) - \epsilon/2\right\}\right) \geq 1/2,$$

and combining this result with (1) shows that $P^{2m}(R) \geq P^m(Q)/2$.

To complete the proof, we show that (2) holds for any $h \in \mathcal{H}$. For a fixed h , notice that $k \text{er}_{\mathbf{w}}(h)$ is a binomial random variable, with expectation $k \text{er}_P(h)$ and variance $\text{er}_P(h)(1 - \text{er}_P(h))k$. Chebyshev's inequality bounds the probability that $\text{er}_{\mathbf{w}}(h) \geq \text{er}_P(h) - \epsilon/2$ by

$$\frac{\text{er}_P(h)(1 - \text{er}_P(h))k}{(\epsilon k/2)^2},$$

which is less than $1/(\epsilon^2 k)$ (using the fact that $x(1-x) \leq 1/4$ for x between 0 and 1). This is at most $1/2$ for $k \geq 2/\epsilon^2$, which implies (2). \square

8.3 Using a group action

Our next task is to seek an upper bound on $P^{n+k}(S_{n,m}^{\epsilon,\gamma}(k))$. We will do this using the next lemma, which is a minor variation on a standard result. Let G_i be the group of all permutations of $\{1, \dots, i\}$ and let σ be a permutation in G_i . For any sequence $\mathbf{a} = (a_1, \dots, a_i)$ we denote by $\sigma \mathbf{a}$ the sequence $(a_{\sigma(1)}, \dots, a_{\sigma(i)})$.

Lemma 19 (see [10], and [3], Lemma 8.3.3) *Let S be a subset of S^{n+k} and let P be any distribution on Z . Let $\mathbf{z}\mathbf{w}$ be a sequence in Z^{n+k} and let G be any subgroup of G_{n+k} . Then,*

$$P^{n+k}(S) \leq \frac{1}{|G|} \max_{\mathbf{z}\mathbf{w} \in Z^{n+k}} N(S, \mathbf{z}\mathbf{w})$$

where $N(S, \mathbf{z}\mathbf{w})$ is defined as

$$N(S, \mathbf{z}\mathbf{w}) = |\{\sigma \in G : \sigma \mathbf{z}\mathbf{w} \in S\}|.$$

FIGURE

Figure 2: Illustration of the effect of the permutations σ_i when applied to $\mathbf{s} = \mathbf{z}\mathbf{w}$.

In the light of lemma 19, our next step is therefore to obtain a bound on $N(S_{n,m}^{\epsilon,\gamma}(k), \mathbf{z}\mathbf{w})$ that is independent of $\mathbf{z}\mathbf{w}$, using a suitable subgroup G . For any $(h_i, h_j) \in \mathcal{H}^2$, define the set,

$$S(h_i, h_j) = \left\{ \mathbf{z}\mathbf{w} \in Z^{n+k} : \text{er}_{\mathbf{z}}^\gamma(h_i) = 0, \text{er}_{\mathbf{z}'}^\gamma(h_j) = 0 \text{ and } \text{er}_{\mathbf{w}}(h_i) \geq \text{er}_{\mathbf{z}_m}^\gamma(h_j) + \frac{\epsilon}{2} \right\}$$

where as usual $\mathbf{z} = \mathbf{z}'\mathbf{z}_m$. We now require that $m \geq 3$, and we choose $k = n' = (n/m)$. We define G as follows. Consider the section $\mathbf{z}_{m-2}\mathbf{z}_{m-1}\mathbf{z}_m\mathbf{w}$ of the sequence $\mathbf{s} = \mathbf{z}\mathbf{w}$ where $\mathbf{s} = (z_1, \dots, z_n, z_{n+1}, \dots, z_{n+n'})$. For any i in the range $(n-3n'+1) \leq i \leq (n-n')$ let σ_i be the permutation of $\{1, \dots, n+n'\}$ having the following properties:

1. $\sigma_i(j) = j$ for $j \neq i$ and $j \neq i + 2n'$.
2. $\sigma_i(i) = i + 2n'$ and $\sigma_i(i + 2n') = i$.

Consequently, if $\mathbf{s} = \mathbf{z}\mathbf{w}$ then each σ_i , when used to obtain $\sigma_i\mathbf{s}$, swaps an element in the $\mathbf{z}_{m-2}\mathbf{z}_{m-1}$ section with a corresponding element in the $\mathbf{z}_m\mathbf{w}$ section. This is illustrated in Figure 2. We define G to be the subgroup of $G_{n+n'}$ generated by the set of permutations $\{\sigma_i\}$ for i in the range defined above. Clearly $|G| = 2^{2n'}$.

8.4 Using covers

We now show how the problem of bounding $N(S, \mathbf{z}\mathbf{w})$ can be reduced to a more manageable problem through the use of covers. Fix $\mathbf{s} = \mathbf{z}\mathbf{w} \in Z^{n+n'}$ and suppose that $\gamma \in (0, 1)$. Suppose that $\mathbf{s} = (z_1, \dots, z_n, z_{n+1}, \dots, z_{n+n'})$ where $z_i = (x_i, y_i) \in Z$. Let \mathcal{C} be a $\gamma/2$ -cover (in the d_∞ sense) of $\mathcal{H}_{|\mathbf{x}|}$, where $\mathbf{x} = (x_1, x_2, \dots, x_{n+n'})$, and suppose that \mathcal{C} has minimum possible cardinality, $\mathcal{N}_\infty(\mathcal{H}, \gamma/2, n+n')$. Then, for each $h \in \mathcal{H}$ there is $\hat{h} \in \mathcal{C}$ such that for $1 \leq i \leq n+n'$, $|h(x_i) - \hat{h}(x_i)| < \gamma/2$. Suppose that \mathbf{s} belongs to $S = S_{n,m}^{\epsilon,\gamma}(n')$. Then, by definition, there are $h, h_m \in \mathcal{H}$ such that

$$\text{er}_{\mathbf{z}'}^\gamma(h_m) = 0, \quad \text{er}_{\mathbf{z}}^\gamma(h) = 0, \quad \text{er}_{\mathbf{w}}(h) > \text{er}_{\mathbf{z}_m}^\gamma(h_m) + \epsilon/2,$$

where, as usual, $\mathbf{z} = \mathbf{z}'\mathbf{z}_m$. If $g \in \mathcal{H}$ satisfies $y_i g(x_i) \geq \gamma$ then, for \hat{g} a member of \mathcal{C} which is $\gamma/2$ -close to g , we have

$$\hat{g}(x_i)y_i \geq g(x_i)y_i - |\hat{g}(x_i) - g(x_i)| > \gamma - \gamma/2 = \gamma/2.$$

Thus the conditions $\text{er}_{\mathbf{z}'}^\gamma(h_m) = 0$ and $\text{er}_{\mathbf{z}}^\gamma(h) = 0$ imply that

$$\text{er}_{\mathbf{z}'}^{\gamma/2}(\hat{h}_m) = \text{er}_{\mathbf{z}}^{\gamma/2}(\hat{h}) = 0,$$

where $\hat{h}_m \in \mathcal{C}$ is $\gamma/2$ -close to h_m , and $\hat{h} \in \mathcal{C}$ is $\gamma/2$ -close to h . Now, if $\text{sgn}(g(x_i)) \neq y_i$ then $g(x_i)y_i < 0$ and $\hat{g}(x_i)y_i < 0 + \gamma/2$. Thus,

$$\text{er}_{\mathbf{w}}^{\gamma/2}(\hat{h}) \geq \text{er}_{\mathbf{w}}(h).$$

Furthermore, if $\hat{g}(x_i)y_i < \gamma/2$ then $g(x_i)y_i < \gamma/2 + \gamma/2 = \gamma$. Therefore,

$$\text{er}_{\mathbf{z}_m}^\gamma(h_m) \geq \text{er}_{\mathbf{z}_m}^{\gamma/2}(\hat{h}_m).$$

These observations show that if $\mathbf{s} = \mathbf{z}'\mathbf{z}_m\mathbf{w} \in S$ then there exist $\hat{h}, \hat{h}_m \in \mathcal{C}$ such that

$$\text{er}_{\mathbf{z}'}^{\gamma/2}(\hat{h}_m) = 0, \quad \text{er}_{\mathbf{z}}^{\gamma/2}(\hat{h}) = 0, \quad \text{er}_{\mathbf{w}}^{\gamma/2}(\hat{h}) > \text{er}_{\mathbf{z}_m}^{\gamma/2}(\hat{h}_m) + \epsilon/2.$$

It follows that

$$N(S, \mathbf{z}\mathbf{w}) \leq N(\hat{S}, \mathbf{z}\mathbf{w}), \tag{3}$$

where

$$\hat{S} = \left\{ \mathbf{z}'\mathbf{z}_m\mathbf{w} : \exists f, g \in \mathcal{C} \text{ s.t. } \text{er}_{\mathbf{z}'}^{\gamma/2}(f) = 0, \quad \text{er}_{\mathbf{z}}^{\gamma/2}(g) = 0, \quad \text{er}_{\mathbf{w}}^{\gamma/2}(g) > \text{er}_{\mathbf{z}_m}^{\gamma/2}(f) + \epsilon/2 \right\}.$$

Now,

$$\begin{aligned} N(\hat{S}, \mathbf{s}) &= \left| \{ \sigma \in G : \sigma\mathbf{s} \in \hat{S} \} \right| \\ &\leq \left| \bigcup_{(f,g) \in \mathcal{C}^2} \{ \sigma \in G : \sigma\mathbf{s} \in \hat{S}(f,g) \} \right| \\ &\leq \sum_{(f,g) \in \mathcal{C}^2} \left| \{ \sigma \in G : \sigma\mathbf{s} \in \hat{S}(f,g) \} \right|, \end{aligned} \tag{4}$$

where

$$\hat{S}(f,g) = \left\{ \mathbf{z}'\mathbf{z}_m\mathbf{w} : \text{er}_{\mathbf{z}'}^{\gamma/2}(f) = 0, \quad \text{er}_{\mathbf{z}}^{\gamma/2}(g) = 0, \quad \text{er}_{\mathbf{w}}^{\gamma/2}(g) > \text{er}_{\mathbf{z}_m}^{\gamma/2}(f) + \epsilon/2 \right\}.$$

We now bound the quantity in the summation of (4) independently of \mathbf{s} and (f,g) .

Lemma 20 *With \mathbf{s}, f, g as above,*

$$\left| \{ \sigma \in G : \sigma\mathbf{s} \in \hat{S}(f,g) \} \right| \leq 2^{2n' - n'\epsilon/2}. \tag{5}$$

Proof Denote by $N(f, g, \mathbf{s})$ the quantity on the left of (5), and suppose it is nonzero (for otherwise, the bound certainly holds). Then there is at least one $\sigma \in G$ such that $\sigma \mathbf{s} = \sigma \mathbf{z} \mathbf{w} \in \hat{S}(f, g)$. Since G is a group, $N(f, g, \mathbf{s})$ is then equal to the number of permutations $\sigma' \in G$ for which $(\sigma' \sigma) \mathbf{s} \in \hat{S}(f, g)$ and consequently we can assume that $\mathbf{s} \in \hat{S}(f, g)$.

Consider any permutation $\sigma \in G$. When used to obtain $\sigma \mathbf{z} \mathbf{w}$ this permutation has the effect of swapping pairs of elements, using combinations of the individual permutations σ_i described above. Whether $\sigma \mathbf{z} \mathbf{w} \in \hat{S}(f, g)$ depends on which pairs are swapped. If σ swaps any pair including an (x_i, y_i) in \mathbf{z}_m for which $f(x_i)y_i < \gamma/2$, then $\sigma \mathbf{z} \mathbf{w} \notin \hat{S}(f, g)$. Similarly, if σ swaps any pair including an (x_i, y_i) in \mathbf{w} such that $g(x_i)y_i < \gamma/2$ then $\sigma \mathbf{z} \mathbf{w} \notin \hat{S}(f, g)$. Consequently,

$$N(f, g, \mathbf{s}) \leq 2^{n_1+n_2}$$

where $n_1 = n' - n' \text{er}_{\mathbf{z}_m}^{\gamma/2}(f)$ is the number of entries (x_i, y_i) of \mathbf{z}_m for which $f(x_i)y_i \geq \gamma/2$, and $n_2 = n' - n' \text{er}_{\mathbf{w}}^{\gamma/2}(g)$ is the number of (x_i, y_i) in \mathbf{w} such that $g(x_i)y_i \geq \gamma/2$. Now, because $\mathbf{s} \in \hat{S}(f, g)$, we have

$$\text{er}_{\mathbf{w}}^{\gamma/2}(g) > \text{er}_{\mathbf{z}_m}^{\gamma/2}(f) + \epsilon/2,$$

so

$$n' - n_2 > n - n_1 + n'\epsilon/2;$$

hence $n_1 - n_2 > n'\epsilon/2$ and

$$n_1 + n_2 = 2n_1 - (n_1 - n_2) < 2n_1 - n'\epsilon/2 \leq 2n' - n'\epsilon/2.$$

It follows, as required, that

$$N(f, g, \mathbf{s}) \leq 2^{2n' - n'\epsilon/2}.$$

□

8.5 Proof of the main theorems

Theorem 6 follows on piecing together the results we have obtained. We have, by (4) and (5),

$$\begin{aligned} N(\hat{S}, \mathbf{s}) &\leq \sum_{(f,g) \in \mathcal{C}^2} 2^{2n' - n'\epsilon/2} \\ &\leq |\mathcal{C}|^2 2^{2n' - n'\epsilon/2} \\ &= (\mathcal{N}_\infty(\mathcal{H}, \gamma/2, n + n'))^2 2^{2n' - n'\epsilon/2}. \end{aligned}$$

By Lemma 18, Lemma 19, and (3), we then have, for $k = n' = n/m \geq 2/\epsilon^2$,

$$\begin{aligned}
P^n(Q_{n,m}^{\epsilon,\gamma}) &\leq 2P^{n+k}(S_{n,m}^{\epsilon,\gamma}(k)) \\
&\leq \frac{2}{|G|} \max_{\mathbf{s} \in Z^{n+k}} N(S, \mathbf{s}) \\
&\leq \frac{2}{|G|} \max_{\mathbf{s} \in Z^{n+k}} N(\hat{S}, \mathbf{s}) \\
&\leq \frac{2}{2^{2n'}} (\mathcal{N}_\infty(\mathcal{H}, \gamma/2, n + n'))^2 2^{2n' - n'\epsilon/2} \\
&= 2 (\mathcal{N}_\infty(\mathcal{H}, \gamma/2, n + n'))^2 2^{-n'\epsilon/2} \\
&= 2 (\mathcal{N}_\infty(\mathcal{H}, \gamma/2, n + n/m))^2 2^{-n\epsilon/2m},
\end{aligned}$$

the required bound.

Theorem 7 follows immediately on using Lemma 17.

To obtain Theorem 8 and Theorem 9 we use a technique from [8]. We give the proof of Theorem 8, that of Theorem 9 being similar.

Proof of Theorem 8: For $0 < \alpha_1, \alpha_2, \delta < 1$, let $E(\alpha_1, \alpha_2, \delta)$ be the set

$$\{\mathbf{x} \in X^n : \exists h_m, h \in \mathcal{H} \text{ s.t. } \text{er}_{\mathbf{z} \setminus \mathbf{z}_m}^{\alpha_2}(h_m) = \text{er}_{\mathbf{z}}^{\alpha_2}(h) = 0, \text{ and } \text{er}_P(h) \geq \text{er}_{\mathbf{x}_m}^{\alpha_2}(h_m) + \epsilon(n, \delta, \alpha_1)\},$$

where

$$\epsilon(n, m, \delta, \alpha_1) = \max \left\{ \sqrt{\frac{2m}{n}}, \frac{2m}{n} \log \left(\frac{2}{\delta} \right) + \frac{4m}{n} \log (\mathcal{N}_\infty(\mathcal{H}, \alpha_1/2, n + n/m)) \right\}.$$

Theorem 6 asserts that $P^n(E(\alpha, \alpha, \delta)) < \delta$ for all α . Furthermore, as is easily checked, if $0 < \alpha_1 \leq \alpha \leq \alpha_2 < 1$ and $\delta_1 \leq \delta$, then $E(\alpha_1, \alpha_2, \delta_1) \subseteq E(\alpha, \alpha, \delta)$. We may argue, as in [8], that

$$\begin{aligned}
P^n \left(\bigcup_{\gamma \in (0,1)} E(\gamma/2, \gamma, \delta\gamma/2) \right) &\leq P^n \left(\bigcup_{i=0}^{\infty} \bigcup_{\gamma \in (2^{-(i+1)}, 2^{-i}]} E(\gamma/2, \gamma, \delta\gamma/2) \right) \\
&\leq P^n \left(\bigcup_{i=0}^{\infty} E(2^{-(i+1)}, 2^{-(i+1)}, \delta 2^{-(i+1)}) \right) \\
&\leq \sum_{i=0}^{\infty} \delta 2^{-(i+1)} = \delta,
\end{aligned}$$

and the result follows on observing that the conclusion of Theorem 8 is precisely that

$$P^n \left(\bigcup_{\gamma \in (0,1)} E(\gamma/2, \gamma, \delta\gamma/2) \right) \leq \delta.$$

9 Conclusions and further work

In this paper we have presented a new cross-validation technique, applicable when real-valued functions are being used for binary classification. We have derived results in probability theory in order to provide bounds on the estimation error of the new technique. As discussed earlier in the paper, the probability results are in a sense stronger than one needs to analyse the cross-validation technique for particular ‘sensible’ learning algorithms, and it is an open question as to whether better estimation error bounds can be found using different methods.

One specific aspect of our bounds which it might be possible to improve concerns the $\sqrt{2m/n}$ term appearing in the estimation error bounds. This arises from the requirement in the proof of Lemma 18 that $k \geq 2/\epsilon^2$. We have been unable to replace this by a condition of the form $k \geq O(1/\epsilon)$, but if this were done then the square-root term could be dropped from the bounds. Perhaps one approach (used for the γ -resubstitution error in [8]) is to develop some result on *relative* deviation of error along the lines of that obtained by Vapnik [25] (see also [7]).

The new holdout and cross-validation techniques we have developed require that each training example be classified correctly, and with a margin, by the relevant function(s). It might be worth allowing a small number of misclassifications. To analyse this in the same vein, we would have to derive probability results similar to those considered here but in which (for the holdout case, for instance), the condition $\text{er}_{\mathbf{z} \setminus \mathbf{z}_m}^\gamma(h_m) = \text{er}_{\mathbf{z}}^\gamma(h) = 0$ is replaced by $\text{er}_{\mathbf{z} \setminus \mathbf{z}_m}^\gamma(h_m), \text{er}_{\mathbf{z}}^\gamma(h) \leq \eta$ for some small η .

Quite apart from theoretical analysis of the performance of the new technique, experimental work remains to be done to test its performance on real data, and this is something we plan to do.

Acknowledgements

Part of this work was carried out while the authors were guests of the Isaac Newton Institute for Mathematical Sciences, Cambridge University. We thank the organisers of the Neural Networks and Machine Learning programme. We are grateful to David Hand for helpful discussion.

References

- [1] N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale-sensitive dimensions, uniform convergence, and learnability. In *Proceedings of the Symposium on Foundations of Computer Science*. IEEE Press, 1993.
- [2] M. Anthony. Probabilistic ‘generalization’ of functions and dimension-based uniform convergence results. To appear, *Statistics and Computing*.
- [3] M. Anthony and N. Biggs. *Computational Learning Theory: An Introduction*. Cambridge Tracts in Theoretical Computer Science (30). Cambridge University Press, Cambridge, UK, 1992.
- [4] M. Anthony and P. Bartlett. *Theory of Learning in Neural Networks* (working title). To be published by Cambridge University Press.
- [5] M. Anthony and P. Bartlett. Function learning from interpolation. Extended abstract in Proceedings EuroCOLT’95, Springer-Verlag, 1995: 211–221.
- [6] N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale-sensitive dimensions, uniform convergence, and learnability. In *Proceedings of the Symposium on Foundations of Computer Science*. IEEE Press, 1993.
- [7] M. Anthony and J. Shawe-Taylor. A result of Vapnik with applications. *Discrete Applied Mathematics*, 47:207–217, 1994.
- [8] P. Bartlett. The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network. Report, Department of Systems Engineering, Australian National University, May 1997.
- [9] P. Bartlett. For valid generalisation, the size of the weights is more important than the size of the network. In *Advances in Neural Information Processing Systems, 9*. Morgan Kaufmann, 1996.
- [10] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *J. ACM*, 36(4):929–965, 1989.
- [11] B. Cheng and D.M. Titterton. Neural networks: a review from a statistical perspective. *Statistical Science*, 9(1):2–54, 1994.
- [12] D. Cohn and G. Tesauro. How tight are the Vapnik-Chervonenkis bounds? *Neural Computation*, 4(2):249–269, 1992.
- [13] R.O. Duda and P.E. Hart. *Pattern Classification and Scene Analysis*. John Wiley, 1973.
- [14] D. Haussler, Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Inform. Comput.*, 100(1):78–150, Sept. 1992.

- [15] S.B. Holden. Cross-validation and the PAC learning model. Research Note RN/96/64, Department of Computer Science, University College London, December 1996.
- [16] S.B. Holden and M. Niranjana. On the practical applicability of VC dimension bounds. *Neural Computation*, 7(6):1265–1288, 1995.
- [17] M. J. Kearns, R. E. Schapire, and L. M. Sellie. Toward efficient agnostic learning. In *Proc. 5th Annu. Workshop on Comput. Learning Theory*, pages 341–352. ACM Press, New York, NY, 1992. **Journal version?**
- [18] M.J. Kearns and D. Ron. Algorithmic stability and sanity-check bounds for leave-one-out cross-validation. In *Proc. 10th Annu. Workshop on Comput. Learning Theory*, pages 152–162. ACM Press, New York, NY, 1997.
- [19] M. J. Kearns and R. E. Schapire. Efficient distribution-free learning of probabilistic concepts. In *Proc. of the 31st Symposium on the Foundations of Comp. Sci.*, pages 382–391. IEEE Computer Society Press, Los Alamitos, CA, 1990. **Journal version?**
- [20] D. Pollard. *Convergence of Stochastic Processes*. Springer-Verlag, 1984.
- [21] N. Sauer. On the density of families of sets. *Journal of Combinatorial Theory (A)*, 13:145–147, 1972.
- [22] J. Shawe-Taylor, P. Bartlett, R.C. Williamson, M. Anthony. Structural risk minimisation over data-dependent hierarchies. To appear, *IEEE Transactions on Information Theory*.
- [23] G.T. Toussaint. Bibliography on estimation of misclassification. *IEEE Transactions on Information Theory*, 20(4):472–479, 1974.
- [24] L.G. Valiant A theory of the learnable. *Commun. ACM*, 27(11):1134–1142, Nov. 1984.
- [25] V. N. Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer-Verlag, New York, 1982.
- [26] V. N. Vapnik and A. Y. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probab. and its Applications*, 16(2):264–280, 1971.