

1

Pattern classification and learning theory

Gábor Lugosi

1.1 A binary classification problem

Pattern recognition (or *classification* or *discrimination*) is about guessing or predicting the unknown class of an observation. An *observation* is a collection of numerical measurements, represented by a d -dimensional vector x . The unknown nature of the observation is called a *class*. It is denoted by y and takes values in the set $\{0, 1\}$. (For simplicity, we restrict our attention to binary classification.) In pattern recognition, one creates a function $g(x) : \mathcal{R}^d \rightarrow \{0, 1\}$ which represents one's guess of y given x . The mapping g is called a *classifier*. A classifier errs on x if $g(x) \neq y$.

To model the learning problem, we introduce a probabilistic setting, and let (X, Y) be an $\mathcal{R}^d \times \{0, 1\}$ -valued random pair.

The random pair (X, Y) may be described in a variety of ways: for example, it is defined by the pair (μ, η) , where μ is the probability measure for X and η is the regression of Y on X . More precisely, for a Borel-measurable set $A \subseteq \mathcal{R}^d$,

$$\mu(A) = \mathbb{P}\{X \in A\},$$

and for any $x \in \mathcal{R}^d$,

$$\eta(x) = \mathbb{P}\{Y = 1 | X = x\} = \mathbb{E}\{Y | X = x\}.$$

Thus, $\eta(x)$ is the conditional probability that Y is 1 given $X = x$. The distribution of (X, Y) is determined by (μ, η) . The function η is called the *a posteriori probability*.

Any function $g : \mathcal{R}^d \rightarrow \{0, 1\}$ defines a *classifier*. An *error* occurs if $g(X) \neq Y$, and the *probability of error* for a classifier g is

$$L(g) = \mathbb{P}\{g(X) \neq Y\}.$$

The Bayes classifier given by

$$g^*(x) = \begin{cases} 1 & \text{if } \eta(x) > 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

minimizes the probability of error:

Theorem 1.1. *For any classifier $g : \mathcal{R}^d \rightarrow \{0, 1\}$,*

$$\mathbb{P}\{g^*(X) \neq Y\} \leq \mathbb{P}\{g(X) \neq Y\}.$$

PROOF. Given $X = x$, the conditional probability of error of any decision g may be expressed as

$$\begin{aligned} & \mathbb{P}\{g(X) \neq Y | X = x\} \\ &= 1 - \mathbb{P}\{Y = g(X) | X = x\} \\ &= 1 - (\mathbb{P}\{Y = 1, g(X) = 1 | X = x\} + \mathbb{P}\{Y = 0, g(X) = 0 | X = x\}) \\ &= 1 - (\mathbb{I}_{\{g(x)=1\}} \mathbb{P}\{Y = 1 | X = x\} + \mathbb{I}_{\{g(x)=0\}} \mathbb{P}\{Y = 0 | X = x\}) \\ &= 1 - (\mathbb{I}_{\{g(x)=1\}} \eta(x) + \mathbb{I}_{\{g(x)=0\}} (1 - \eta(x))), \end{aligned}$$

where \mathbb{I}_A denotes the indicator of the set A . Thus, for every $x \in \mathcal{R}^d$,

$$\begin{aligned} & \mathbb{P}\{g(X) \neq Y | X = x\} - \mathbb{P}\{g^*(X) \neq Y | X = x\} \\ &= \eta(x) (\mathbb{I}_{\{g^*(x)=1\}} - \mathbb{I}_{\{g(x)=1\}}) + (1 - \eta(x)) (\mathbb{I}_{\{g^*(x)=0\}} - \mathbb{I}_{\{g(x)=0\}}) \\ &= (2\eta(x) - 1) (\mathbb{I}_{\{g^*(x)=1\}} - \mathbb{I}_{\{g(x)=1\}}) \\ &\geq 0 \end{aligned}$$

by the definition of g^* . The statement now follows by integrating both sides with respect to $\mu(dx)$. \square

L^* is called the Bayes probability of error, Bayes error, or Bayes risk. The proof above reveals that

$$L(g) = 1 - \mathbb{E}\{\mathbb{I}_{\{g(X)=1\}} \eta(X) + \mathbb{I}_{\{g(X)=0\}} (1 - \eta(X))\},$$

and in particular,

$$L^* = 1 - \mathbb{E}\{\mathbb{I}_{\{\eta(X)>1/2\}} \eta(X) + \mathbb{I}_{\{\eta(X)\leq 1/2\}} (1 - \eta(X))\} = \mathbb{E}\min(\eta(X), 1 - \eta(X)).$$

Note that g^* depends upon the distribution of (X, Y) . If this distribution is known, g^* may be computed. Most often, the distribution of (X, Y) is unknown, so that g^* is unknown too.

In our model, we have access to a data base of pairs (X_i, Y_i) , $1 \leq i \leq n$, observed in the past. We assume that $(X_1, Y_1), \dots, (X_n, Y_n)$, the *data*, is a sequence of independent identically distributed (*i.i.d.*) random pairs with the same distribution as that of (X, Y) .

A classifier is constructed on the basis of $X_1, Y_1, \dots, X_n, Y_n$ and is denoted by g_n : Y is guessed by $g_n(X; X_1, Y_1, \dots, X_n, Y_n)$. The process of constructing g_n is called *learning*, *supervised learning*, or *learning with a teacher*. The performance of g_n is measured by the conditional *probability of error*

$$L_n = L(g_n) = \mathbb{P}\{g_n(X; X_1, Y_1, \dots, X_n, Y_n) \neq Y | X_1, Y_1, \dots, X_n, Y_n\}.$$

This is a random variable because it depends upon the data. So, L_n averages over the distribution of (X, Y) , but the data is held fixed. Even though averaging over the data as well is unnatural, since in a given application, one has to live with the data at hand, the number $\mathbb{E}L_n = \mathbb{P}\{g_n(X) \neq Y\}$ which indicates the quality on an average data sequence, provides useful information, especially if the random variable L_n is concentrated around its mean with high probability.

1.2 Empirical risk minimization

Assume that a class \mathcal{C} of classifiers $g : \mathcal{R}^d \rightarrow \{0, 1\}$ is given and our task is to find one with a small probability of error. In the lack of the knowledge of the underlying distribution, one has to resort to using the data to estimate the probabilities of error for the classifiers in \mathcal{C} . It is tempting to pick a classifier from \mathcal{C} that minimizes an estimate of the probability of error over the class. The most natural choice to estimate the probability of error $L(g) = \mathbb{P}\{g(X) \neq Y\}$ is the error count

$$\hat{L}_n(g) = \frac{1}{n} \sum_{j=1}^n \mathbb{I}_{\{g(X_j) \neq Y_j\}}.$$

$\hat{L}_n(g)$ is called the *empirical error* of the classifier g .

A good method should pick a classifier with a probability of error that is close to the minimal probability of error in the class. Intuitively, if we can estimate the error probability for the classifiers in \mathcal{C} *uniformly* well, then the classification function that minimizes the estimated probability of error is likely to have a probability of error that is close to the best in the class.

Denote by g_n^* the classifier that minimizes the estimated probability of error over the class:

$$\hat{L}_n(g_n^*) \leq \hat{L}_n(g) \quad \text{for all } g \in \mathcal{C}.$$

Then for the probability of error

$$L(g_n^*) = \mathbb{P}\{g_n^*(X) \neq Y | D_n\}$$

of the selected rule we have:

Lemma 1.1.

$$\begin{aligned} L(g_n^*) - \inf_{g \in \mathcal{C}} L(g) &\leq 2 \sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)|, \\ |\hat{L}_n(g_n^*) - L(g_n^*)| &\leq \sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)|. \end{aligned}$$

PROOF.

$$\begin{aligned} L(g_n^*) - \inf_{g \in \mathcal{C}} L(g) &= L(g_n^*) - \hat{L}_n(g_n^*) + \hat{L}_n(g_n^*) - \inf_{g \in \mathcal{C}} L(g) \\ &\leq L(g_n^*) - \hat{L}_n(g_n^*) + \sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)| \\ &\leq 2 \sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)|. \end{aligned}$$

The second inequality is trivially true. \square

We see that upper bounds for $\sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)|$ provide us with upper bounds for two things simultaneously:

- (1) An upper bound for the suboptimality of g_n^* within \mathcal{C} , that is, a bound for $L(g_n^*) - \inf_{g \in \mathcal{C}} L(g)$.
- (2) An upper bound for the error $|\hat{L}_n(g_n^*) - L(g_n^*)|$ committed when $\hat{L}_n(g_n^*)$ is used to estimate the probability of error $L(g_n^*)$ of the selected rule.

It is particularly useful to know that even though $\hat{L}_n(g_n^*)$ is usually optimistically biased, it is within given bounds of the unknown probability of error with g_n^* , and that no other test sample is needed to estimate this probability of error. Whenever our bounds indicate that we are close to the optimum in \mathcal{C} , we must at the same time have a good estimate of the probability of error, and vice versa.

The random variable $n\hat{L}_n(g)$ is binomially distributed with parameters n and $L(g)$. Thus, to obtain bounds for the success of empirical error minimization, we need to study uniform deviations of binomial random variables from their means. In the next two sections we summarize the basics of the underlying theory.

1.3 Concentration inequalities

1.3.1 Hoeffding's inequality

The simplest inequality to bound the difference between a random variable and its expected value is *Markov's inequality*: for any nonnegative random variable X , and $t > 0$,

$$\mathbb{P}\{X \geq t\} \leq \frac{\mathbb{E}X}{t}.$$

From this, we deduce *Chebyshev's inequality*: if X is an arbitrary random variable and $t > 0$, then

$$\mathbb{P}\{|X - \mathbb{E}X| \geq t\} = \mathbb{P}\{|X - \mathbb{E}X|^2 \geq t^2\} \leq \frac{\mathbb{E}\{|X - \mathbb{E}X|^2\}}{t^2} = \frac{\text{Var}\{X\}}{t^2}.$$

As an example, we derive inequalities for $\mathbb{P}\{S_n - \mathbb{E}S_n \geq t\}$ with $S_n = \sum_{i=1}^n X_i$, where X_1, \dots, X_n are independent real-valued random variables. Chebyshev's inequality and independence immediately gives us

$$\mathbb{P}\{|S_n - \mathbb{E}S_n| \geq t\} \leq \frac{\text{Var}\{S_n\}}{t^2} = \frac{\sum_{i=1}^n \text{Var}\{X_i\}}{t^2}.$$

The meaning of this is perhaps better seen if we assume that the X_i 's are i.i.d. Bernoulli(p) random variables (i.e., $\mathbb{P}\{X_i = 1\} = 1 - \mathbb{P}\{X_i = 0\} = p$), and normalize:

$$\mathbb{P}\left\{\left|\frac{1}{n} \sum_{i=1}^n X_i - p\right| \geq \epsilon\right\} \leq \frac{p(1-p)}{n\epsilon^2}.$$

To illustrate the weakness of this bound, let $\Phi(y) = \int_{-\infty}^y e^{-t^2/2}/\sqrt{2\pi} dt$ be the normal distribution function. The central limit theorem states that

$$\mathbb{P}\left\{\sqrt{\frac{n}{p(1-p)}} \left(\frac{1}{n} \sum_{i=1}^n X_i - p\right) \geq y\right\} \rightarrow 1 - \Phi(y) \leq \frac{1}{\sqrt{2\pi}} \frac{e^{-y^2/2}}{y},$$

from which we would expect something like

$$\mathbb{P}\left\{\frac{1}{n} \sum_{i=1}^n X_i - p \geq \epsilon\right\} \approx e^{-n\epsilon^2/(2p(1-p))}.$$

Clearly, Chebyshev's inequality is off mark. An improvement may be obtained by *Chernoff's bounding method*. By Markov's inequality, if s is an arbitrary positive number, then for any random variable X , and any $t > 0$,

$$\mathbb{P}\{X \geq t\} = \mathbb{P}\{e^{sX} \geq e^{st}\} \leq \frac{\mathbb{E}e^{sX}}{e^{st}}.$$

In Chernoff's method, we find an $s > 0$ that minimizes the upper bound or makes the upper bound small. In the case of a sum of independent random variables,

$$\begin{aligned} \mathbb{P}\{S_n - \mathbb{E}S_n \geq t\} &\leq e^{-st} \mathbb{E}\left\{\exp\left(s \sum_{i=1}^n (X_i - \mathbb{E}X_i)\right)\right\} \\ &= e^{-st} \prod_{i=1}^n \mathbb{E}\left\{e^{s(X_i - \mathbb{E}X_i)}\right\} \quad (\text{by independence}). \end{aligned}$$

Now the problem of finding tight bounds comes down to finding a good upper bound for the moment generating function of the random variables $X_i - \mathbb{E}X_i$. There are many ways

of doing this. For bounded random variables perhaps the most elegant version is due to Hoeffding (1963):

Lemma 1.2. *Let X be a random variable with $\mathbb{E}X = 0$, $a \leq X \leq b$. Then for $s > 0$,*

$$\mathbb{E}\{e^{sX}\} \leq e^{s^2(b-a)^2/8}.$$

PROOF. Note that by convexity of the exponential function

$$e^{sx} \leq \frac{x-a}{b-a}e^{sb} + \frac{b-x}{b-a}e^{sa} \quad \text{for } a \leq x \leq b.$$

Exploiting $\mathbb{E}X = 0$, and introducing the notation $p = -a/(b-a)$ we get

$$\begin{aligned} \mathbb{E}e^{sX} &\leq \frac{b}{b-a}e^{sa} - \frac{a}{b-a}e^{sb} \\ &= (1-p + pe^{s(b-a)})e^{-ps(b-a)} \\ &\stackrel{\text{def}}{=} e^{\phi(u)}, \end{aligned}$$

where $u = s(b-a)$, and $\phi(u) = -pu + \log(1-p + pe^u)$. But by straightforward calculation it is easy to see that the derivative of ϕ is

$$\phi'(u) = -p + \frac{p}{p + (1-p)e^{-u}},$$

therefore $\phi(0) = \phi'(0) = 0$. Moreover,

$$\phi''(u) = \frac{p(1-p)e^{-u}}{(p + (1-p)e^{-u})^2} \leq \frac{1}{4}.$$

Thus, by Taylor series expansion with remainder, for some $\theta \in [0, u]$,

$$\phi(u) = \phi(0) + u\phi'(0) + \frac{u^2}{2}\phi''(\theta) \leq \frac{u^2}{8} = \frac{s^2(b-a)^2}{8}.$$

Now we may directly plug this lemma into the bound obtained by Chernoff's method:

$$\begin{aligned} &\mathbb{P}\{S_n - \mathbb{E}S_n \geq \epsilon\} \\ &\leq e^{-s\epsilon} \prod_{i=1}^n \mathbb{E}\{e^{s(X_i - \mathbb{E}X_i)}\} \\ &\leq e^{-s\epsilon} \prod_{i=1}^n e^{s^2(b_i - a_i)^2/8} \quad (\text{by Lemma 1.2}) \\ &= e^{-s\epsilon} e^{s^2 \sum_{i=1}^n (b_i - a_i)^2/8} \\ &= e^{-2\epsilon^2 / \sum_{i=1}^n (b_i - a_i)^2} \quad (\text{by choosing } s = 4\epsilon / \sum_{i=1}^n (b_i - a_i)^2). \end{aligned}$$

The result we have just derived is generally known as *Hoeffding's inequality*. For binomial random variables it was proved by Chernoff (1952) and Okamoto (1952). Summarizing, we have:

Theorem 1.2. (HOEFFDING'S INEQUALITY). *Let X_1, \dots, X_n be independent bounded random variables such that X_i falls in the interval $[a_i, b_i]$ with probability one. Denote their sum by $S_n = \sum_{i=1}^n X_i$. Then for any $\epsilon > 0$ we have*

$$\mathbb{P}\{S_n - \mathbb{E}S_n \geq \epsilon\} \leq e^{-2\epsilon^2 / \sum_{i=1}^n (b_i - a_i)^2}$$

and

$$\mathbb{P}\{S_n - \mathbb{E}S_n \leq -\epsilon\} \leq e^{-2\epsilon^2 / \sum_{i=1}^n (b_i - a_i)^2}.$$

If we specialize this to the binomial distribution, that is, when the X_i 's are i.i.d. Bernoulli(p), we get

$$\mathbb{P}\{S_n/n - p \geq \epsilon\} \leq e^{-2n\epsilon^2},$$

which is just the kind of inequality we hoped for.

We may combine this inequality with that of Lemma 1.1 to bound the performance of empirical risk minimization in the special case when the class \mathcal{C} contains finitely many classifiers:

Theorem 1.3. *Assume that the cardinality of \mathcal{C} is bounded by N . Then we have for all $\epsilon > 0$,*

$$\mathbb{P}\left\{\sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)| > \epsilon\right\} \leq 2Ne^{-2n\epsilon^2}.$$

An important feature of the result above is that it is completely distribution free. The actual distribution of the data does not play a role at all in the upper bound.

To have an idea about the size of the error, one may be interested in the expected maximal deviation

$$\mathbb{E} \sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)|.$$

The inequality above may be used to derive such an upper bound by observing that for any nonnegative random variable X ,

$$\mathbb{E}X = \int_0^\infty \mathbb{P}\{X \geq t\} dt.$$

Sharper bounds result by combining Lemma 1.2 with the following simple result:

Lemma 1.3. *Let $\sigma > 0$, $n \geq 2$, and let Y_1, \dots, Y_n be real-valued random variables such that for all $s > 0$ and $1 \leq i \leq n$, $\mathbb{E}\{e^{sY_i}\} \leq e^{s^2\sigma^2/2}$. Then*

$$\mathbb{E}\left\{\max_{i \leq n} Y_i\right\} \leq \sigma\sqrt{2\ln n}.$$

If, in addition, $\mathbb{E}\{e^{s(-Y_i)}\} \leq e^{s^2\sigma^2/2}$ for every $s > 0$ and $1 \leq i \leq n$, then for any $n \geq 1$,

$$\mathbb{E}\left\{\max_{i \leq n}|Y_i|\right\} \leq \sigma\sqrt{2\ln(2n)}.$$

PROOF. By Jensen's inequality, for all $s > 0$,

$$e^{s\mathbb{E}\{\max_{i \leq n} Y_i\}} \leq \mathbb{E}\{e^{s\max_{i \leq n} Y_i}\} = \mathbb{E}\left\{\max_{i \leq n} e^{sY_i}\right\} \leq \sum_{i=1}^n \mathbb{E}\{e^{sY_i}\} \leq ne^{s^2\sigma^2/2}.$$

Thus,

$$\mathbb{E}\left\{\max_{i \leq n} Y_i\right\} \leq \frac{\ln n}{s} + \frac{s\sigma^2}{2},$$

and taking $s = \sqrt{2\ln n/\sigma^2}$ yields the first inequality. Finally, note that $\max_{i \leq n} |Y_i| = \max(Y_1, -Y_1, \dots, Y_n, -Y_n)$ and apply the first inequality to prove the second. \square

Now we obtain

$$\mathbb{E}\sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)| \leq \sqrt{\frac{\ln(2N)}{2n}}.$$

1.3.2 Other inequalities for sums

Here we summarize some other useful inequalities for the deviations of sums of independent random variables from their means.

Theorem 1.4. BENNETT'S INEQUALITY. Let X_1, \dots, X_n be independent real-valued random variables with zero mean, and assume that $|X_i| \leq c$ with probability one. Let $\sigma^2 = \frac{1}{n} \sum_{i=1}^n \text{Var}\{X_i\}$. Then that for any $t > 0$,

$$\mathbb{P}\{S_n > t\} \leq \exp\left(-\frac{n\sigma^2}{c^2} h\left(\frac{ct}{n\sigma^2}\right)\right),$$

where the function h is defined by $h(u) = (1+u)\log(1+u) - u$ for $u \geq 0$.

SKETCH OF PROOF. We use Chernoff's method as in the proof of Hoeffding's inequality. Write

$$\mathbb{E}\{e^{sX_i}\} = 1 + s\mathbb{E}\{X_i\} + \sum_{r=2}^{\infty} \frac{s^r \mathbb{E}\{X_i^r\}}{r!} = 1 + s^2 \text{Var}\{X_i\} F_i \leq e^{s^2 \text{Var}\{X_i\} F_i}$$

with $F_i = \sum_{r=2}^{\infty} s^{r-2} \mathbb{E}\{X_i^r\} / (r! \text{Var}\{X_i\})$. We may use the boundedness of the X_i 's to show that $\mathbb{E}\{X_i^r\} \leq c^{r-2} \text{Var}\{X_i\}$, which implies $F_i \leq (e^{sc} - 1 - sc) / (sc)^2$. Choose the s which minimizes the obtained upper bound for the tail probability. \square

Theorem 1.5. BERNSTEIN'S INEQUALITY. *Under the conditions of the previous exercise, for any $t > 0$,*

$$\mathbb{P}\{S_n > t\} \leq \exp\left(-\frac{t^2}{2n\sigma^2 + 2ct/3}\right).$$

PROOF. The result follows from Bennett's inequality and the inequality $h(u) \geq u^2/(2+2u/3)$, $u \geq 0$. \square

Theorem 1.6. Let X_1, \dots, X_n be independent random variables, taking their values from $[0, 1]$. If $m = \mathbb{E}S_n$, then for any $m \leq t \leq n$,

$$\mathbb{P}\{S_n \geq t\} \leq \left(\frac{m}{t}\right)^t \left(\frac{n-m}{n-t}\right)^{n-t}.$$

Also,

$$\mathbb{P}\{S_n \geq t\} \leq \left(\frac{m}{t}\right)^t e^{t-m},$$

and for all $\epsilon > 0$,

$$\mathbb{P}\{S_n \geq m(1 + \epsilon)\} \leq e^{-mh(\epsilon)},$$

where h is the function defined in the previous theorem. Finally,

$$\mathbb{P}\{S_n \leq m(1 - \epsilon)\} \leq e^{-m\epsilon^2/2}.$$

1.3.3 The bounded difference inequality

In this section we give some powerful extensions of concentration inequalities for sums to general functions of independent random variables.

Let A be some set, and let $g : A^n \rightarrow \mathcal{R}$ be some measurable function of n variables. We derive inequalities for the difference between $g(X_1, \dots, X_n)$ and its expected value when X_1, \dots, X_n are arbitrary independent random variables taking values in A . Sometimes we will write g instead of $g(X_1, \dots, X_n)$ whenever it does not cause any confusion.

We recall the elementary fact that if X and Y are arbitrary bounded random variables, then $\mathbb{E}\{XY\} = \mathbb{E}\{\mathbb{E}\{XY|Y\}\} = \mathbb{E}\{Y\mathbb{E}\{X|Y\}\}$.

The first result of this section is an improvement of an inequality of Efron and Stein (1981) proved by Steele (1986). We have learnt the short proof given here from Stéphane Boucheron.

Theorem 1.7. EFRON-STEIN INEQUALITY. *If X'_1, \dots, X'_n form an independent copy of X_1, \dots, X_n , then*

$$\mathbf{Var}(g(X_1, \dots, X_n)) \leq \frac{1}{2} \sum_{i=1}^n \mathbb{E}\{(g(X_1, \dots, X_n) - g(X_1, \dots, X'_i, \dots, X_n))^2\}$$

PROOF. Introduce the notation $V = g - \mathbb{E}g$, and define

$$V_i = \mathbb{E}\{g|X_1, \dots, X_i\} - \mathbb{E}\{g|X_1, \dots, X_{i-1}\}, \quad i = 1, \dots, n.$$

Clearly, $V = \sum_{i=1}^n V_i$. Then

$$\begin{aligned} \mathbf{Var}(g) &= \mathbb{E} \left\{ \left(\sum_{i=1}^n V_i \right)^2 \right\} \\ &= \mathbb{E} \sum_{i=1}^n V_i^2 + 2\mathbb{E} \sum_{i>j} V_i V_j \\ &= \mathbb{E} \sum_{i=1}^n V_i^2, \end{aligned}$$

since, for any $i > j$,

$$\mathbb{E}V_i V_j = \mathbb{E}\mathbb{E}\{V_i V_j | X_1, \dots, X_j\} = \mathbb{E}\{V_j \mathbb{E}\{V_i | X_1, \dots, X_j\}\} = 0.$$

To bound $\mathbb{E}V_i^2$, note that, by Jensen's inequality,

$$\begin{aligned} V_i^2 &= (\mathbb{E}\{g|X_1, \dots, X_i\} - \mathbb{E}\{g|X_1, \dots, X_{i-1}\})^2 \\ &= \left(\mathbb{E} \left[\mathbb{E}\{g|X_1, \dots, X_n\} - \mathbb{E}\{g|X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\} \middle| X_1, \dots, X_i \right] \right)^2 \\ &\leq \mathbb{E} \left[(\mathbb{E}\{g|X_1, \dots, X_n\} - \mathbb{E}\{g|X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\})^2 \middle| X_1, \dots, X_i \right], \end{aligned}$$

and therefore

$$\begin{aligned} \mathbb{E}V_i^2 &\leq \mathbb{E} \left[(g - \mathbb{E}\{g|X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\})^2 \right] \\ &= \frac{1}{2} \mathbb{E} \left[(g(X_1, \dots, X_n) - g(X_1, \dots, X'_i, \dots, X_n))^2 \right], \end{aligned}$$

where at the last step we used (conditionally) the elementary fact that if X and Y are independent and identically distributed random variables, then $\mathbf{Var}(X) = (1/2)\mathbb{E}\{(X - Y)^2\}$. \square

Assume that a function $g : A^n \rightarrow \mathcal{R}$ satisfies the *bounded difference assumption*

$$\sup_{\substack{x_1, \dots, x_n, \\ x'_i \in A}} |g(x_1, \dots, x_n) - g(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c_i, \quad 1 \leq i \leq n.$$

In other words, we assume that if we change the i -th variable of g while keeping all the others fixed, then the value of the function does not change by more than c_i . Then the Efron-Stein inequality implies that

$$\mathbf{Var}(g) \leq \frac{1}{2} \sum_{i=1}^n c_i^2.$$

For such functions it is possible to prove the following exponential tail inequality, a powerful extension of Hoeffding's inequality.

Theorem 1.8. THE BOUNDED DIFFERENCE INEQUALITY. *Under the bounded difference assumption above, for all $t > 0$,*

$$\mathbb{P}\{g(X_1, \dots, X_n) - \mathbb{E}g(X_1, \dots, X_n) \geq t\} \leq e^{-2t^2/\sum_{i=1}^n c_i^2},$$

and

$$\mathbb{P}\{\mathbb{E}g(X_1, \dots, X_n) - g(X_1, \dots, X_n) \geq t\} \leq e^{-2t^2/\sum_{i=1}^n c_i^2}.$$

McDiarmid (1989) proved this inequality using martingale techniques, which we reproduce here. The proof of Theorem 1.8 uses the following straightforward extension of Lemma 1.2:

Lemma 1.4. *Let V and Z be random variables such that $\mathbb{E}\{V|Z\} = 0$ with probability one, and for some function h and constant $c \geq 0$*

$$h(Z) \leq V \leq h(Z) + c.$$

Then for all $s > 0$

$$\mathbb{E}\{e^{sV}|Z\} \leq e^{s^2c^2/8}.$$

PROOF OF THEOREM 1.8. Just like in the proof of Theorem 1.7, introduce the notation $V = g - \mathbb{E}g$, and define

$$V_i = \mathbb{E}\{g|X_1, \dots, X_i\} - \mathbb{E}\{g|X_1, \dots, X_{i-1}\}, \quad i = 1, \dots, n.$$

Then $V = \sum_{i=1}^n V_i$. Also introduce the random variables

$$H_i(X_1, \dots, X_i) = \mathbb{E}\{g(X_1, \dots, X_n)|X_1, \dots, X_i\}.$$

Then, denoting the distribution of X_i by F_i for $i = 1, \dots, n$,

$$V_i = H_i(X_1, \dots, X_i) - \int H_i(X_1, \dots, X_{i-1}, x)F_i(dx).$$

Define the random variables

$$W_i = \sup_u \left(H_i(X_1, \dots, X_{i-1}, u) - \int H_i(X_1, \dots, X_{i-1}, x)F_i(dx) \right),$$

and

$$Z_i = \inf_v \left(H_i(X_1, \dots, X_{i-1}, v) - \int H_i(X_1, \dots, X_{i-1}, x)F_i(dx) \right).$$

Clearly, $Z_i \leq V_i \leq W_i$ with probability one, and also

$$W_i - Z_i = \sup_u \sup_v (H_i(X_1, \dots, X_{i-1}, u) - H_i(X_1, \dots, X_{i-1}, v)) \leq c_i ,$$

by the bounded difference assumption. Therefore, we may apply the lemma above to obtain, for all $i = 1, \dots, n$,

$$\mathbb{E}\{e^{sV_i} | X_1, \dots, X_{i-1}\} \leq e^{s^2 c_i^2 / 8}.$$

Finally, by Chernoff's bound, for any $s > 0$,

$$\begin{aligned} & \mathbb{P}\{g - \mathbb{E}g \geq t\} \\ & \leq \frac{\mathbb{E}\{e^{s\sum_{i=1}^n V_i}\}}{e^{st}} = \frac{\mathbb{E}\left\{e^{s\sum_{i=1}^{n-1} V_i} \mathbb{E}\{e^{sV_n} | X_1, \dots, X_{n-1}\}\right\}}{e^{st}} \\ & \leq e^{s^2 c_n^2 / 8} \frac{\mathbb{E}\{e^{s\sum_{i=1}^{n-1} V_i}\}}{e^{st}} \\ & \leq e^{-st} e^{s^2 \sum_{i=1}^n c_i^2 / 8} \quad (\text{by repeating the same argument } n \text{ times}). \end{aligned}$$

Choosing $s = 4t / \sum_{i=1}^n c_i^2$ proves the first inequality. The proof of the second inequality is similar. \square

An important application of the bounded difference inequality shows that if \mathcal{C} is any class of classifiers of form $g : \mathcal{R}^d \rightarrow \{0, 1\}$, then

$$\mathbb{P}\left\{\left|\sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)| - \mathbb{E}\sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)|\right| > \epsilon\right\} \leq 2e^{-2n\epsilon^2}.$$

Indeed, if we view $\sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)|$ as a function of the n independent random pairs (X_i, Y_i) , $i = 1, \dots, n$, then we immediately see that the bounded difference assumption is satisfied with $c_i = 1/n$, and Theorem 1.8 immediately implies the statement.

The interesting fact is that regardless of the size of its expected value, the random variable $\sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)|$ is sharply concentrated around its mean with very large probability. In the next section we study the expected value.

1.4 Vapnik-Chervonenkis theory

1.4.1 The Vapnik-Chervonenkis inequality

Recall from Section 1.3.1 that for any finite class \mathcal{C} of classifiers, and for all $\epsilon > 0$,

$$\mathbb{P}\left\{\sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)| > \epsilon\right\} \leq 2Ne^{-2n\epsilon^2},$$

and

$$\mathbb{E} \sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)| \leq \sqrt{\frac{\ln(2N)}{2n}}.$$

These simple bounds may be useless if the cardinality N of the class is very large, or infinite. The purpose of this section is to introduce a theory to handle such cases.

Let X_1, \dots, X_n be i.i.d. random variables taking values in \mathcal{R}^d with common distribution

$$\mu(A) = \mathbb{P}\{X_1 \in A\} \quad (A \subset \mathcal{R}^d).$$

Define the empirical distribution

$$\mu_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{[X_i \in A]} \quad (A \subset \mathcal{R}^d).$$

Consider a class \mathcal{A} of subsets of \mathcal{R}^d . Our main concern here is the behavior of the random variable $\sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)|$. We saw in the previous chapter that a simple consequence of the bounded difference inequality is that

$$\mathbb{P} \left\{ \left| \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| - \mathbb{E} \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right| > t \right\} \leq 2e^{-2nt^2}$$

for any n and $t > 0$. This shows that for any class \mathcal{A} , the maximal deviation is sharply concentrated around its mean. In the rest of this chapter we derive inequalities for the expected value, in terms of certain combinatorial quantities related to \mathcal{A} . The first such quantity is the vc *shatter coefficient*, defined by

$$\mathbb{S}_{\mathcal{A}}(n) = \max_{x_1, \dots, x_n \in \mathcal{R}^d} |\{ \{x_1, \dots, x_n\} \cap A; A \in \mathcal{A} \}|.$$

Thus, $\mathbb{S}_{\mathcal{A}}(n)$ is the maximal number of different subsets of a set of n points which can be obtained by intersecting it with elements of \mathcal{A} . The main theorem is the following version of a classical result of Vapnik and Chervonenkis:

Theorem 1.9. VAPNIK-CHERVENENKIS INEQUALITY.

$$\mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right\} \leq 2\sqrt{\frac{\log 2\mathbb{S}_{\mathcal{A}}(n)}{n}}.$$

PROOF. Introduce X'_1, \dots, X'_n , an independent copy of X_1, \dots, X_n . Also, define n i.i.d. sign variables $\sigma_1, \dots, \sigma_n$ such that $\mathbb{P}\{\sigma_1 = -1\} = \mathbb{P}\{\sigma_1 = 1\} = 1/2$, independent of

$X_1, X'_1, \dots, X_n, X'_n$. Then, denoting $\mu'_n(A) = (1/n) \sum_{i=1}^n \mathbb{I}_{[X'_i \in A]}$, we may write

$$\begin{aligned}
& \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right\} \\
&= \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mathbb{E}\{\mu_n(A) - \mu'_n(A) | X_1, \dots, X_n\}| \right\} \\
&\leq \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \mathbb{E} \left\{ |\mu_n(A) - \mu'_n(A)| \mid X_1, \dots, X_n \right\} \right\} \\
&\quad (\text{by Jensen's inequality}) \\
&\leq \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu'_n(A)| \right\} \\
&\quad (\text{since } \sup \mathbb{E}(\cdot) \leq \mathbb{E} \sup(\cdot)) \\
&= \frac{1}{n} \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbb{I}_{[X_i \in A]} - \mathbb{I}_{[X'_i \in A]}) \right| \right\} \\
&\quad (\text{because } X_1, X'_1, \dots, X_n, X'_n \text{ are i.i.d.}) \\
&= \frac{1}{n} \mathbb{E} \left\{ \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbb{I}_{[x_i \in A]} - \mathbb{I}_{[x'_i \in A]}) \right| \mid X_1, X'_1, \dots, X_n, X'_n \right\} \right\}.
\end{aligned}$$

Now because of the independence of the σ_i 's of the rest of the variables, we may fix the values of $X_1 = x_1, X'_1 = x'_1, \dots, X_n = x_n, X'_n = x'_n$, and investigate

$$\mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbb{I}_{[x_i \in A]} - \mathbb{I}_{[x'_i \in A]}) \right| \right\}.$$

Denote by $\widehat{\mathcal{A}} \subset \mathcal{A}$ a collection of sets such that any two sets in $\widehat{\mathcal{A}}$ have different intersections with the set $\{x_1, x'_1, \dots, x_n, x'_n\}$, and every possible intersection is represented once. Thus, $|\widehat{\mathcal{A}}| \leq \mathbb{S}_{\mathcal{A}}(2n)$, and

$$\mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbb{I}_{[x_i \in A]} - \mathbb{I}_{[x'_i \in A]}) \right| \right\} = \mathbb{E} \left\{ \max_{A \in \widehat{\mathcal{A}}} \left| \sum_{i=1}^n \sigma_i (\mathbb{I}_{[x_i \in A]} - \mathbb{I}_{[x'_i \in A]}) \right| \right\}.$$

Observing that each $\sigma_i (\mathbb{I}_{[x_i \in A]} - \mathbb{I}_{[x'_i \in A]})$ has zero mean and takes values in $[-1, 1]$, we obtain from Lemma 1.2 that for any $s > 0$,

$$\mathbb{E} e^{s \sum_{i=1}^n \sigma_i (\mathbb{I}_{[x_i \in A]} - \mathbb{I}_{[x'_i \in A]})} = \prod_{i=1}^n \mathbb{E} e^{s \sigma_i (\mathbb{I}_{[x_i \in A]} - \mathbb{I}_{[x'_i \in A]})} \leq e^{ns^2/2}.$$

Since the distribution of $\sigma_i (\mathbb{I}_{[x_i \in A]} - \mathbb{I}_{[x'_i \in A]})$ is symmetric, Lemma 1.3 immediately implies that

$$\mathbb{E} \left\{ \max_{A \in \widehat{\mathcal{A}}} \left| \sum_{i=1}^n \sigma_i (\mathbb{I}_{[x_i \in A]} - \mathbb{I}_{[x'_i \in A]}) \right| \right\} \leq \sqrt{2n \log 2\mathbb{S}_{\mathcal{A}}(2n)}.$$

Conclude by observing that $\mathbb{S}_{\mathcal{A}}(2n) \leq \mathbb{S}_{\mathcal{A}}(n)^2$. \square

Remark. The original form of the Vapnik-Chervonenkis inequality is

$$\mathbb{P} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| > t \right\} \leq 4\mathbb{S}_{\mathcal{A}}(2n)e^{-nt^2/8}.$$

A combination of Theorem 1.9 with the concentration inequality for the supremum quickly yields an inequality of a similar form.

The main virtue of the Vapnik-Chervonenkis inequality is that it converts the problem of uniform deviations of empirical averages into a combinatorial problem. Investigating the behavior of $\mathbb{S}_{\mathcal{A}}(n)$ is the key to the understanding of the behavior of the maximal deviations. Classes for which $\mathbb{S}_{\mathcal{A}}(n)$ grows at a subexponential rate with n are manageable in the sense that $\mathbb{E}\{\sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)|\}$ converges to zero. More importantly, explicit upper bounds for $\mathbb{S}_{\mathcal{A}}(n)$ provide nonasymptotic distribution-free bounds for the expected maximal deviation (and also for the tail probabilities). Section 1.4.3 is devoted to some key combinatorial results related to shatter coefficients.

We close this section by a refinement of Theorem 1.9 due to Massart (2000). The bound below substantially improves the bound of Theorem 1.9 whenever $\sup_{A \in \mathcal{A}} \mu(A)(1 - \mu(A))$ is very small.

Theorem 1.10. Let $\Sigma = \sup_{A \in \mathcal{A}} \sqrt{\mu(A)(1 - \mu(A))}$. Then

$$\mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right\} \leq \frac{16 \log 2\mathbb{S}_{\mathcal{A}}(2n)}{n} + \sqrt{\frac{32\Sigma^2 \log 2\mathbb{S}_{\mathcal{A}}(2n)}{n}}.$$

PROOF. From the proof of Theorem 1.9, we have

$$\begin{aligned} & \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right\} \\ & \leq \frac{1}{n} \mathbb{E} \left\{ \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbb{I}_{[X_i \in A]} - \mathbb{I}_{[X'_i \in A]}) \right| \middle| X_1, X'_1, \dots, X_n, X'_n \right\} \right\}. \end{aligned}$$

By Hoeffding's inequality, for each set A ,

$$\mathbb{E} \left\{ e^{s \sum_{i=1}^n \sigma_i (\mathbb{I}_{[X_i \in A]} - \mathbb{I}_{[X'_i \in A]})} \middle| X_1, X'_1, \dots, X_n, X'_n \right\} \leq e^{s^2 \sum_{i=1}^n (\mathbb{I}_{[X_i \in A]} - \mathbb{I}_{[X'_i \in A]})^2 / 2},$$

so by Lemma 1.3 we obtain

$$\mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right\} \leq \frac{1}{n} \mathbb{E} \sup_{A \in \mathcal{A}} \sqrt{\sum_{i=1}^n (\mathbb{I}_{[X_i \in A]} - \mathbb{I}_{[X'_i \in A]})^2} \sqrt{2 \log 2\mathbb{S}_{\mathcal{A}}(2n)}.$$

To bound the right-hand side, note that

$$\begin{aligned}
& \mathbb{E} \sup_{A \in \mathcal{A}} \sqrt{\sum_{i=1}^n (\mathbb{I}_{[X_i \in A]} - \mathbb{I}_{[X'_i \in A]})^2} \\
& \leq \sqrt{\mathbb{E} \sup_{A \in \mathcal{A}} \sum_{i=1}^n (\mathbb{I}_{[X_i \in A]} - \mathbb{I}_{[X'_i \in A]})^2} \\
& \leq \sqrt{\mathbb{E} \sup_{A \in \mathcal{A}} \sum_{i=1}^n ((\mathbb{I}_{[X_i \in A]} - \mu(A)) + (\mu(A) - \mathbb{I}_{[X'_i \in A]}))^2} \\
& \leq \sqrt{4\mathbb{E} \sup_{A \in \mathcal{A}} \sum_{i=1}^n (\mathbb{I}_{[X_i \in A]} - \mu(A))^2} \\
& = 2\sqrt{\mathbb{E} \sup_{A \in \mathcal{A}} \sum_{i=1}^n [(\mathbb{I}_{[X_i \in A]} - \mu(A))(1 - \mu(A)) + \mu(A)(\mu(A) - \mathbb{I}_{[X_i \in A]}) + \mu(A)(1 - \mu(A))]} \\
& \leq 2\sqrt{n\Sigma^2} + 2\sqrt{\mathbb{E} \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n (\mathbb{I}_{[X_i \in A]} - \mu(A)) \right|} \\
& = 2\sqrt{n\Sigma^2} + 2\sqrt{n\mathbb{E} \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)|}.
\end{aligned}$$

Summarizing, if we denote $\mathbb{E} \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| = M$, we have obtained

$$M \leq \sqrt{\frac{\log 2\mathbb{S}_{\mathcal{A}}(2n)}{2n}} (\Sigma + \sqrt{M}).$$

This is a quadratic inequality for \sqrt{M} , whose solution is just the statement of the theorem.

□

1.4.2 Inequalities for relative deviations

In this section we summarize some important improvements of the basic Vapnik-Chervonenkis inequality. The basic result is the following pair of inequalities, due to Vapnik and Chervonenkis (1974). The proof sketched here is due to Anthony and Shawe-Taylor (1993).

Theorem 1.11. *For every $\epsilon > 0$,*

$$\mathbb{P} \left\{ \sup_{A \in \mathcal{A}} \frac{\mu(A) - \mu_n(A)}{\sqrt{\mu(A)}} > \epsilon \right\} \leq 4\mathbb{S}_{\mathcal{A}}(2n)e^{-n\epsilon^2/4}$$

and

$$\mathbb{P} \left\{ \sup_{A \in \mathcal{A}} \frac{\mu_n(A) - \mu(A)}{\sqrt{\mu_n(A)}} > \epsilon \right\} \leq 4S_{\mathcal{A}}(2n)e^{-n\epsilon^2/4}.$$

SKETCH OF PROOF. The main steps of the proof are as follows:

1. Symmetrization.

$$\mathbb{P} \left\{ \sup_{A \in \mathcal{A}} \frac{\mu(A) - \mu_n(A)}{\sqrt{\mu(A)}} > \epsilon \right\} \leq 2\mathbb{P} \left\{ \sup_{A \in \mathcal{A}} \frac{\mu'_n(A) - \mu_n(A)}{\sqrt{(1/2)(\mu'_n(A) + \mu_n(A))}} > \epsilon \right\}.$$

2. Randomization, conditioning.

$$\begin{aligned} & \mathbb{P} \left\{ \sup_{A \in \mathcal{A}} \frac{\mu'_n(A) - \mu_n(A)}{\sqrt{(1/2)(\mu'_n(A) + \mu_n(A))}} > \epsilon \right\} \\ &= \mathbb{E} \left\{ \mathbb{P} \left\{ \sup_{A \in \mathcal{A}} \frac{(1/n) \sum_{i=1}^n \sigma_i (\mathbb{I}_{X_i \in A} - \mathbb{I}_{X_i \notin A})}{\sqrt{(1/2)(\mu'_n(A) + \mu_n(A))}} > \epsilon \mid X_1, X'_1, \dots, X_n, X'_n \right\} \right\}. \end{aligned}$$

3. Tail bound. Use the union bound and Hoeffding's inequality to bound the conditional probability inside. \square

Using the bounds above, we may derive other interesting inequalities. The first inequalities are due to Pollard (1995) and Haussler (1992).

COROLLARY 1.1. *For all $t \in (0, 1)$ and $s > 0$,*

$$\mathbb{P} \left\{ \sup_{A \in \mathcal{A}} \frac{\mu(A) - \mu_n(A)}{\mu(A) + \mu_n(A) + s/2} > t \right\} \leq 4S_{\mathcal{A}}(2n)e^{-nst^2/4}$$

and

$$\mathbb{P} \left\{ \sup_{A \in \mathcal{A}} \frac{\mu_n(A) - \mu(A)}{\mu(A) + \mu_n(A) + s/2} > t \right\} \leq 4S_{\mathcal{A}}(2n)e^{-nst^2/4}.$$

SKETCH OF PROOF. Take $\alpha > 0$. Considering the cases $\mu(A) < (\alpha + 1)^2 \epsilon^2 \alpha^{-2}$ and $\mu(A) \geq (\alpha + 1)^2 \epsilon^2 \alpha^{-2}$ separately, it is easy to show that $\mu(A) - \mu_n(A) \leq \epsilon \sqrt{\mu(A)}$ implies that $\mu(A) \leq (1 + \alpha)\mu_n(A) + \epsilon^2(1 + \alpha)/\alpha$. Then choosing $\alpha = 2t/(1-t)$ and $\epsilon^2 = st^2/(1-t^2)$ we easily prove that the first inequality in Theorem 1.11 implies the first inequality. The second inequality follows similarly from the second inequality of Theorem 1.11. \square

Finally, we point out another corollary of Theorem 1.11 which has interesting applications in statistical learning theory:

COROLLARY 1.2.

$$\mathbb{P}\{\exists A \in \mathcal{A} : \mu(A) > \epsilon \text{ and } \mu_n(A) \leq (1-t)\mu(A)\} \leq 4\mathbb{S}_{\mathcal{A}}(2n)e^{-n\epsilon t^2/4}.$$

In particular, setting $t = 1$,

$$\mathbb{P}\{\exists A \in \mathcal{A} : \mu(A) > \epsilon \text{ and } \mu_n(A) = 0\} \leq 4\mathbb{S}_{\mathcal{A}}(2n)e^{-n\epsilon/4}.$$

1.4.3 Shatter coefficients

Consider a class \mathcal{A} of subsets of \mathcal{R}^d , and let $x_1, \dots, x_n \in \mathcal{R}^d$ be arbitrary points. Recall from the previous section that properties of the finite set $\mathcal{A}(x_1^n) \subset \{0,1\}^n$ defined by

$$\begin{aligned} \mathcal{A}(x_1^n) = \{b = (b_1, \dots, b_n) \in \{0,1\}^n : \\ b_i = \mathbb{I}_{[x_i \in A]}, i = 1, \dots, n \text{ for some } A \in \mathcal{A}\} \end{aligned}$$

play an essential role in bounding uniform deviations of the empirical measure. In particular, the maximal cardinality of $\mathcal{A}(x_1^n)$

$$\mathbb{S}_{\mathcal{A}}(n) = \max_{x_1, \dots, x_n \in \mathcal{R}^d} |\mathcal{A}(x_1^n)|$$

(i.e., the shatter coefficient) provides simple bounds via the Vapnik-Chervonenkis inequality. We begin with some elementary properties of the shatter coefficient.

Theorem 1.12. *Let \mathcal{A} and \mathcal{B} be classes of subsets of \mathcal{R}^d , and let $n, m \geq 1$ be integers. Then*

- (1) $\mathbb{S}_{\mathcal{A}}(n+m) \leq \mathbb{S}_{\mathcal{A}}(n)\mathbb{S}_{\mathcal{A}}(m)$;
- (2) If $\mathcal{C} = \mathcal{A} \cup \mathcal{B}$, then $\mathbb{S}_{\mathcal{C}}(n) \leq \mathbb{S}_{\mathcal{A}}(n) + \mathbb{S}_{\mathcal{B}}(n)$;
- (3) If $\mathcal{C} = \{C = A^c : A \in \mathcal{A}\}$, then $\mathbb{S}_{\mathcal{C}}(n) = \mathbb{S}_{\mathcal{A}}(n)$;
- (4) If $\mathcal{C} = \{C = A \cap B : A \in \mathcal{A} \text{ and } B \in \mathcal{B}\}$, then $\mathbb{S}_{\mathcal{C}}(n) \leq \mathbb{S}_{\mathcal{A}}(n)\mathbb{S}_{\mathcal{B}}(n)$;
- (5) If $\mathcal{C} = \{C = A \cup B : A \in \mathcal{A} \text{ and } B \in \mathcal{B}\}$, then $\mathbb{S}_{\mathcal{C}}(n) \leq \mathbb{S}_{\mathcal{A}}(n)\mathbb{S}_{\mathcal{B}}(n)$;
- (6) If $\mathcal{C} = \{C = A \times B : A \in \mathcal{A} \text{ and } B \in \mathcal{B}\}$, then $\mathbb{S}_{\mathcal{C}}(n) \leq \mathbb{S}_{\mathcal{A}}(n)\mathbb{S}_{\mathcal{B}}(n)$.

PROOF. Parts (1), (2), (3), and (6) are immediate from the definition. To show (4), fix x_1, \dots, x_n , let $N = |\mathcal{A}(x_1^n)| \leq \mathbb{S}_{\mathcal{A}}(n)$, and denote by A_1, A_2, \dots, A_N the different sets of the form $\{x_1, \dots, x_n\} \cap A$ for some $A \in \mathcal{A}$. For all $1 \leq i \leq N$, sets in \mathcal{B} pick at most $\mathbb{S}_{\mathcal{B}}(|A_i|) \leq \mathbb{S}_{\mathcal{B}}(n)$ different subsets of A_i . Thus,

$$|\mathcal{A}(x_1^n)| \leq \sum_{i=1}^N \mathbb{S}_{\mathcal{B}}(|A_i|) \leq \mathbb{S}_{\mathcal{A}}(n)\mathbb{S}_{\mathcal{B}}(n).$$

(5) follows from (4) and (3). \square

The VC *dimension* V of a class \mathcal{A} of sets is defined as the largest integer n such that

$$\mathbb{S}_{\mathcal{A}}(n) = 2^n.$$

If $\mathbb{S}_{\mathcal{A}}(n) = 2^n$ for all n , then we say that $V = \infty$. Clearly, if $\mathbb{S}_{\mathcal{A}}(n) < 2^n$ for some n , then for all $m > n$, $\mathbb{S}_{\mathcal{A}}(m) < 2^m$, and therefore the VC dimension is always well-defined. If $|\mathcal{A}(x_1^n)| = 2^n$ for some points x_1, \dots, x_n , then we say that \mathcal{A} *shatters* the set $x_1^n = \{x_1, \dots, x_n\}$. As the next basic result shows, the VC dimension provides a useful bound for the shatter coefficient of a class.

Theorem 1.13. SAUER'S LEMMA. *Let \mathcal{A} be a class of sets with VC dimension $V < \infty$. Then for all n ,*

$$\mathbb{S}_{\mathcal{A}}(n) \leq \sum_{i=0}^V \binom{n}{i}.$$

PROOF. Fix x_1, \dots, x_n , such that $|\mathcal{A}(x_1^n)| = \mathbb{S}_{\mathcal{A}}(n)$. Denote $B_0 = \mathcal{A}(x_1^n) \in \{0, 1\}^n$. We say that a set $B \subset \{0, 1\}^n$ *shatters* a set $S = \{s_1, \dots, s_m\} \subset \{1, 2, \dots, n\}$ if the restriction of B to the components s_1, \dots, s_m is the full m -dimensional binary hypercube, that is,

$$\{(b_{s_1}, \dots, b_{s_m}) : b = (b_1, \dots, b_n) \in B\} = \{0, 1\}^m.$$

It suffices to show that the cardinality of any set $B_0 \subset \{0, 1\}^n$ that cannot shatter any set of size $m > V$, is at most $\sum_{i=0}^V \binom{n}{i}$. This is done by transforming B_0 into a set B_n with $|B_n| = |B_0|$ such that any set shattered by B_n is also shattered by B_0 . Moreover, it will be easy to see that $|B_n| \leq \sum_{i=0}^V \binom{n}{i}$.

For every vector $b = (b_1, \dots, b_n) \in B_0$, if $b_1 = 1$, then flip the first component of b to zero unless $(0, b_2, \dots, b_n) \in B_0$. If $b_1 = 0$, then keep the vector unchanged. The set of vectors B_1 obtained this way obviously has the same cardinality as that of B_0 . Moreover, if B_1 shatters a set $S = \{s_1, s_2, \dots, s_m\} \subset \{1, \dots, n\}$, then B_0 also shatters S . This is trivial if $1 \notin S$. If $1 \in S$, then we may assume without loss of generality that $s_1 = 1$. The fact that B_1 shatters S implies that for any $v \in \{0, 1\}^{m-1}$ there exists a $b \in B_1$ such that $b_1 = 1$ and $(b_{s_2}, \dots, b_{s_m}) = v$. By the construction of B_1 this is only possible if for any $u \in \{0, 1\}^m$ there exists a $b' \in B_0$ such that $(b'_{s_1}, \dots, b'_{s_m}) = u$. This means that B_0 also shatters S .

Now starting from B_1 , execute the same transformation, but now by flipping the second component of each vector, if necessary. Again, the cardinality of the obtained set B_2 remains unchanged, and any set shattered by B_2 is also shattered by B_1 (and therefore also by B_0). Repeat the transformation for all components, arriving at the set B_n . Clearly, B_n cannot shatter sets of cardinality larger than V , since otherwise B_0 would shatter sets of the same

size. On the other hand, it is easy to see that B_n is such that for every $b \in B_n$, all vectors of form $c = (c_1, \dots, c_n)$ with $c_i \in \{b_i, 0\}$ for $1 \leq i \leq n$, are also in B_n . Then B_n is a subset of a set of form

$$T = \{b \in \{0, 1\}^n : b_i = 0 \text{ if } v_i = 0\},$$

where $v = (v_1, \dots, v_n)$ is a fixed vector containing at most V 1's. This implies that

$$\mathbb{S}_{\mathcal{A}}(n) = |B_0| = |B_n| \leq |T| = \sum_{i=0}^V \binom{n}{i},$$

concluding the proof. \square

The following corollary makes the meaning of Sauer's lemma more transparent:

COROLLARY 1.3. *Let \mathcal{A} be a class of sets with VC dimension $V < \infty$. Then for all n ,*

$$\mathbb{S}_{\mathcal{A}}(n) \leq (n+1)^V,$$

and for all $n \geq V$,

$$\mathbb{S}_{\mathcal{A}}(n) \leq \left(\frac{ne}{V}\right)^V.$$

PROOF. By the binomial theorem,

$$(n+1)^V = \sum_{i=0}^V n^i \binom{V}{i} = \sum_{i=0}^V \frac{n^i V!}{i!(V-i)!} \geq \sum_{i=0}^V \frac{n^i}{i!} \geq \sum_{i=0}^V \binom{n}{i}.$$

On the other hand, if $V/n \leq 1$, then

$$\left(\frac{V}{n}\right)^V \sum_{i=0}^V \binom{n}{i} \leq \sum_{i=0}^V \left(\frac{V}{n}\right)^i \binom{n}{i} \leq \sum_{i=0}^n \left(\frac{V}{n}\right)^i \binom{n}{i} = \left(1 + \frac{V}{n}\right)^n \leq e^V,$$

where again we used the binomial theorem. \square

Recalling the Vapnik-Chervonenkis inequality, we see that if \mathcal{A} is any class of sets with VC dimension V , then

$$\mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right\} \leq 2 \sqrt{\frac{V \log(n+1) + \log 2}{n}},$$

that is, whenever \mathcal{A} has a finite VC dimension, the expected largest deviation over \mathcal{A} converges to zero at a rate $O(\sqrt{\log n/n})$.

Next we calculate the VC dimension of some simple classes.

Lemma 1.5. *If \mathcal{A} is the class of all rectangles in \mathbb{R}^d , then $V = 2d$.*

PROOF. To see that there are $2d$ points that can be shattered by \mathcal{A} , just consider the $2d$ vectors with $d - 1$ zero components, and one non-zero component which is either 1 or -1 . On the other hand, for any given set of $2d + 1$ points we can choose a subset of at most $2d$ points with the property that it contains a point with largest first coordinate, a point with smallest first coordinate, a point with largest second coordinate, and so forth. Clearly, there is no set in \mathcal{A} which contains these points, but not the rest. \square

Lemma 1.6. *Let \mathcal{G} be an m -dimensional vector space of real-valued functions defined on \mathcal{R}^d . The class of sets*

$$\mathcal{A} = \{\{x : g(x) \geq 0\} : g \in \mathcal{G}\}$$

has VC dimension $V \leq m$.

PROOF. It suffices to show that no set of size $m + 1$ can be shattered by sets of the form $\{x : g(x) \geq 0\}$. Fix $m + 1$ arbitrary points x_1, \dots, x_{m+1} , and define the linear mapping $L : \mathcal{G} \rightarrow \mathcal{R}^{m+1}$ as

$$L(g) = (g(x_1), \dots, g(x_{m+1})) .$$

Then the image of \mathcal{G} , $L(\mathcal{G})$, is a linear subspace of \mathcal{R}^{m+1} of dimension not exceeding m . This implies the existence of a nonzero vector $\gamma = (\gamma_1, \dots, \gamma_{m+1}) \in \mathcal{R}^{m+1}$ orthogonal to $L(\mathcal{G})$, that is, for every $g \in \mathcal{G}$,

$$\gamma_1 g(x_1) + \dots + \gamma_{m+1} g(x_{m+1}) = 0 .$$

We may assume that at least one of the γ_i 's is negative. Rearranging this equality so that all terms with nonnegative γ_i stay on the left-hand side, we get

$$\sum_{i: \gamma_i \geq 0} \gamma_i g(x_i) = \sum_{i: \gamma_i < 0} -\gamma_i g(x_i) .$$

Now suppose that there exists a $g \in \mathcal{G}$ such that the set $\{x : g(x) \geq 0\}$ picks exactly the x_i 's on the left-hand side. Then all terms on the left-hand side are nonnegative, while the terms on the right-hand side must be negative, which is a contradiction, so x_1, \dots, x_{m+1} cannot be shattered, which implies the statement. \square

Generalizing a result of Schläffli (1950), Cover (1965) showed that if \mathcal{G} is defined as the linear space of functions spanned by functions $\psi_1, \dots, \psi_m : \mathcal{R}^d \rightarrow \mathcal{R}$, and the vectors $\Psi(x_i) = (\psi_1(x_i), \dots, \psi_m(x_i))$, $i = 1, 2, \dots, n$ are linearly independent, then for the class of sets $\mathcal{A} = \{\{x : g(x) \geq 0\} : g \in \mathcal{G}\}$ we have

$$|\mathcal{A}(x_1^n)| = 2 \sum_{i=0}^{m-1} \binom{n-1}{i} ,$$

which often gives a slightly sharper estimate than Sauer's lemma. The proof is left as an exercise. Now we may immediately deduce the following:

COROLLARY 1.4. (1) If \mathcal{A} is the class of all linear halfspaces, that is, subsets of \mathbb{R}^d of the form $\{x : a^T x \geq b\}$, where $a \in \mathbb{R}^d, b \in \mathbb{R}$ take all possible values, then $V \leq d + 1$.

(2) If \mathcal{A} is the class of all closed balls in \mathbb{R}^d , that is, sets of the form

$$\left\{x = (x^{(1)}, \dots, x^{(d)}) : \sum_{i=1}^d |x^{(i)} - a_i|^2 \leq b\right\}, \quad a_1, \dots, a_d, b \in \mathbb{R},$$

then $V \leq d + 2$.

(3) If \mathcal{A} is the class of all ellipsoids in \mathbb{R}^d , that is, sets of form $\{x : x^T \Sigma^{-1} x \leq 1\}$, where Σ is a positive definite symmetric matrix, then $V \leq d(d + 1)/2 + 1$.

Note that the above-mentioned result implies that the VC dimension of the class of all linear halfspaces actually equals $d + 1$. Dudley (1979) proved that in the case of the class of all closed balls the above inequality is not tight, and the VC dimension equals $d + 1$ (see exercise 5).

1.4.4 Applications to empirical risk minimization

In this section we apply the main results of the previous sections to obtain upper bounds for the performance of empirical risk minimization.

Recall the scenario set up in Chapter 2: \mathcal{C} is a class of classifiers containing decision functions of the form $g : \mathbb{R}^d \rightarrow \{0, 1\}$. The data $(X_1, Y_1), \dots, (X_n, Y_n)$ may be used to calculate the empirical error $\hat{L}_n(g)$ for any $g \in \mathcal{C}$. g_n^* denotes a classifier minimizing $\hat{L}_n(g)$ over the class, that is,

$$\hat{L}_n(g_n^*) \leq \hat{L}_n(g) \quad \text{for all } g \in \mathcal{C}.$$

Denote the probability of error of the optimal classifier in the class by $L_{\mathcal{C}}$, that is,

$$L_{\mathcal{C}} = \inf_{g \in \mathcal{C}} L(g).$$

(Here we implicitly assume that the infimum is achieved. This assumption is motivated by convenience in the notation, it is not essential.)

The basic Lemma 1.1 shows that

$$L(g_n^*) - L_{\mathcal{C}} \leq 2 \sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)|.$$

Thus, the quantity of interest is the maximal deviation between empirical probabilities of error and their expectation over the class. Such quantities are estimated by the Vapnik-Chervonenkis inequality. Indeed, the random variable $\sup_{g \in \mathcal{C}} |\hat{L}_n(g) - L(g)|$ is of the form of $\sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)|$, where the role of the class of sets \mathcal{A} is now played by the class of error sets

$$\{(x, y) \in \mathcal{R}^d \times \{0, 1\} : g(x) \neq y\}; \quad g \in \mathcal{C}.$$

Denote the class of these error sets by $\bar{\mathcal{A}}$. Thus, the Vapnik-Chervonenkis inequality immediately bounds the expected maximal deviation in terms of the shatter coefficients (or VC dimension) of the class of error sets.

Instead of error sets, it is more convenient to work with classes of sets of the form

$$\{x \in \mathcal{R}^d : g(x) = 1\}; \quad g \in \mathcal{C}.$$

We denote the class of sets above by \mathcal{A} . The next simple fact shows that the classes $\bar{\mathcal{A}}$ and \mathcal{A} are equivalent from a combinatorial point of view:

Lemma 1.7. *For every n we have $\mathbb{S}_{\bar{\mathcal{A}}}(n) = \mathbb{S}_{\mathcal{A}}(n)$, and therefore the corresponding VC dimensions are also equal: $V_{\bar{\mathcal{A}}} = V_{\mathcal{A}}$.*

PROOF. Let N be a positive integer. We show that for any n pairs from $\mathcal{R}^d \times \{0, 1\}$, if N sets from $\bar{\mathcal{A}}$ pick N different subsets of the n pairs, then there are N corresponding sets in \mathcal{A} that pick N different subsets of n points in \mathcal{R}^d , and vice versa. Fix n pairs $(x_1, 0), \dots, (x_m, 0), (x_{m+1}, 1), \dots, (x_n, 1)$. Note that since ordering does not matter, we may arrange any n pairs in this manner. Assume that for a certain set $A \in \mathcal{A}$, the corresponding set $\bar{A} = A \times \{0\} \cup A^c \times \{1\} \in \bar{\mathcal{A}}$ picks out the pairs $(x_1, 0), \dots, (x_k, 0), (x_{m+1}, 1), \dots, (x_{m+l}, 1)$, that is, the set of these pairs is the intersection of \bar{A} and the n pairs. Again, we can assume without loss of generality that the pairs are ordered in this way. This means that A picks from the set $\{x_1, \dots, x_n\}$ the subset $\{x_1, \dots, x_k, x_{m+l+1}, \dots, x_n\}$, and the two subsets uniquely determine each other. This proves $\mathbb{S}_{\bar{\mathcal{A}}}(n) \leq \mathbb{S}_{\mathcal{A}}(n)$. To prove the other direction, notice that if A picks a subset of k points x_1, \dots, x_k , then the corresponding set $\bar{A} \in \bar{\mathcal{A}}$ picks the pairs with the same indices from $\{(x_1, 0), \dots, (x_k, 0)\}$. Equality of the VC dimensions follows from the equality of the shatter coefficients. \square

From this point on, we will denote the common value of $\mathbb{S}_{\bar{\mathcal{A}}}(n)$ and $\mathbb{S}_{\mathcal{A}}(n)$ by $\mathbb{S}_{\mathcal{C}}(n)$, and refer to it as the n -th shatter coefficient of the class \mathcal{C} . It is simply the maximum number of different ways n points can be classified by classifiers in the class \mathcal{C} . Similarly, $V_{\bar{\mathcal{A}}} = V_{\mathcal{A}}$ will be referred to as the VC dimension of the class \mathcal{C} , and will be denoted by $V_{\mathcal{C}}$.

Now we are prepared to summarize our main performance bound for empirical risk minimization:

COROLLARY 1.5.

$$\mathbb{E}L(g_n^*) - L_C \leq 4\sqrt{\frac{\log 2S_C(n)}{n}} \leq 4\sqrt{\frac{V_C \log(n+1) + \log 2}{n}}$$

Bounds for $\mathbb{P}\{L(g_n^*) - L_C > \epsilon\}$ may now be easily obtained by combining the corollary above with the bounded difference inequality.

The inequality above may be improved in various different ways. In the appendix of this chapter we show that the factor of $\log n$ in the upper bound is unnecessary, it may be replaced by a suitable constant. In practice, however, often the sample size is so small that the inequality above provides smaller numerical values.

On the other hand, the main performance may be improved in another direction. To understand the reason, consider first an extreme situation when $L_C = 0$, that is, there exists a classifier in \mathcal{C} which classifies without error. (This also means that for some $g' \in \mathcal{C}$, $Y = g'(X)$ with probability one, a very restrictive assumption. Nevertheless, the assumption that $L_C = 0$ is common in computational learning theory, see Blumer, Ehrenfeucht, Haussler, and Warmuth (1989). In such a case, clearly $\hat{L}_n(g^*) = 0$, and the second statement of Corollary 1.2 implies that

$$\mathbb{P}\{L(g_n^*) - L_C > \epsilon\} = \mathbb{P}\{L(g_n^*) > \epsilon\} \leq 4S_C(2n)e^{-n\epsilon/4},$$

and therefore

$$\mathbb{E}L(g_n^*) - L_C = \mathbb{E}L(g_n^*) \leq \frac{4\ln(4S_C(2n))}{n}.$$

(The bound on the expected value may be obtained by the following simple bounding argument: assume that for some nonnegative random variable Z , for all $\epsilon > 0$, $\mathbb{P}\{Z > \epsilon\} \leq Ce^{-K\epsilon}$ for some positive constants. Then $\mathbb{E}Z = \int_0^\infty \mathbb{P}\{Z > \epsilon\}d\epsilon \leq u + \int_u^\infty Ce^{-K\epsilon} d\epsilon$ for any $u > 0$. Integrating, and choosing u to minimize the upper bound, we obtain $\mathbb{E}Z \leq \ln C/K$.)

The main point here is that the upper bound obtained in this special case is of smaller order of magnitude than in the general case ($O(V_C \ln n/n)$ as opposed to $O(\sqrt{V_C \ln n/n})$). Intuition suggests that if L_C is nonzero but very small, the general bound of Corollary 1.5 should be improvable. In fact, the argument below shows that it is possible interpolate between the special case $L_C = 0$ and the fully distribution-free bound of Corollary 1.5:

Theorem 1.14.

$$\mathbb{E}L(g_n^*) - L_C \leq \sqrt{\frac{8L_C \ln(5S_C(2n)) + 2}{n}} + \frac{8\ln(10S_C(2n)) + 4}{n}.$$

Also, for every $\epsilon > 0$,

$$\mathbb{P}\{L(g_n^*) - L_C > \epsilon\} \leq 5S_C(2n)e^{-n\epsilon^2/16(L_C + \epsilon)}.$$

PROOF. For any $\epsilon > 0$, if

$$\sup_{g \in \mathcal{C}} \frac{L(g) - \hat{L}_n(g)}{\sqrt{L(g)}} \leq \frac{\epsilon}{\sqrt{L_C + 2\epsilon}},$$

then for each $g \in \mathcal{C}$

$$\hat{L}_n(g) \geq L(g) - \epsilon \sqrt{\frac{L(g)}{L_C + 2\epsilon}}.$$

If, in addition, g is such that $L(g) > L_C + 2\epsilon$, then by the monotonicity of the function $x - c\sqrt{x}$ (for $c > 0$ and $x > c^2/4$),

$$\hat{L}_n(g) \geq L_C + 2\epsilon - \epsilon \sqrt{\frac{L_C + 2\epsilon}{L_C + 2\epsilon}} = L_C + \epsilon.$$

Therefore,

$$\mathbb{P} \left\{ \inf_{g: L(g) > L_C + 2\epsilon} \hat{L}_n(g) < L_C + \epsilon \right\} \leq \mathbb{P} \left\{ \sup_{g \in \mathcal{C}} \frac{L(g) - \hat{L}_n(g)}{\sqrt{L(g)}} > \frac{\epsilon}{\sqrt{L_C + 2\epsilon}} \right\}.$$

But if $L(g_n^*) - L_C > 2\epsilon$, then, denoting by g' a classifier in \mathcal{C} such that $L(g') = L_C$, there exists an $g \in \mathcal{C}$ such that $L(g) > L_C + 2\epsilon$ and $\hat{L}_n(g) \leq \hat{L}_n(g')$. Thus,

$$\begin{aligned} & \mathbb{P}\{L(g_n^*) - L_C > 2\epsilon\} \\ & \leq \mathbb{P} \left\{ \inf_{g: L(g) > L_C + 2\epsilon} \hat{L}_n(g) < \hat{L}_n(g') \right\} \\ & \leq \mathbb{P} \left\{ \inf_{g: L(g) > L_C + 2\epsilon} \hat{L}_n(g) < L_C + \epsilon \right\} + \mathbb{P}\{\hat{L}_n(g') > L_C + \epsilon\} \\ & \leq \mathbb{P} \left\{ \sup_{g \in \mathcal{C}} \frac{L(g) - \hat{L}_n(g)}{\sqrt{L(g)}} > \frac{\epsilon}{\sqrt{L_C + 2\epsilon}} \right\} + \mathbb{P}\{\hat{L}_n(g') - L_C > \epsilon\}. \end{aligned}$$

Bounding the last two probabilities by Theorem 1.11 and Bernstein's inequality, respectively, we obtain the probability bound of the statement.

The upper bound for the expected value may now be derived by some straightforward calculations which we sketch here: let $u \leq L_C$ be a positive number. Then, using the tail

inequality obtained above,

$$\begin{aligned}
& \mathbb{E}L(g_n^*) - L_C \\
&= \int_0^\infty \mathbb{P}\{L(g_n^*) - L_C > \epsilon\} d\epsilon \\
&\leq u + \int_u^\infty 5S_C(2n) \max\left(e^{-n\epsilon^2/8L_C}, e^{-n\epsilon/8}\right) d\epsilon \\
&\leq \left(u/2 + \int_u^\infty 5S_C(2n) e^{-n\epsilon^2/8L_C} d\epsilon\right) \\
&\quad + \left(u/2 + \int_u^\infty 5S_C(2n) e^{-n\epsilon/8} d\epsilon\right).
\end{aligned}$$

The second term may be bounded as in the argument given for the case $L_C = 0$, while the first term may be calculated similarly, using the additional observation that

$$\begin{aligned}
\int_u^\infty e^{-n\epsilon^2} d\epsilon &\leq \frac{1}{2} \int_u^\infty \left(2 + \frac{1}{n\epsilon^2}\right) e^{-n\epsilon^2} d\epsilon \\
&= \frac{1}{2} \left[\frac{1}{n\epsilon} e^{-n\epsilon^2} \right]_u^\infty.
\end{aligned}$$

The details are omitted. \square

1.4.5 Convex combinations of classifiers

Several important classification methods form a classifier as a convex combination of simple functions. To describe such a situation, consider a class \mathcal{C} of classifiers $g : \mathcal{R}^d \rightarrow \{0, 1\}$. Think of \mathcal{C} as a small class of “base” classifiers such as the class of all linear splits of \mathcal{R}^d . In general we assume that the VC dimension $V_{\mathcal{C}}$ of \mathcal{C} is finite. Define the class \mathcal{F} as the class of functions $f : \mathcal{R}^d \rightarrow [0, 1]$ of the form

$$f(x) = \sum_{j=1}^N w_j g_j(x)$$

where N is any positive integer, w_1, \dots, w_N are nonnegative weights with $\sum_{j=1}^N w_j = 1$, and $g_1, \dots, g_N \in \mathcal{C}$. Thus, \mathcal{F} may be considered as the convex hull of \mathcal{C} . Each function $f \in \mathcal{F}$ defines a classifier g_f , in a natural way, by

$$g_f(x) = \begin{cases} 1 & \text{if } f(x) > 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

A large variety of “boosting” and “bagging” methods, based mostly on the work of Schapire (1990), Freund (1995) and Breiman (1996), construct classifiers as convex combinations

of very simple functions. Typically the class of classifiers defined this way is too large in the sense that it is impossible to obtain meaningful distribution-free upper bounds for $\sup_{f \in \mathcal{F}} (L(g_f) - \hat{L}_n(g_f))$. Indeed, even in the simple case when $d = 1$ and \mathcal{C} is the class of all linear splits of the real line, the class of all g_f is easily seen to have an infinite VC dimension.

Surprisingly, however, meaningful bounds may be obtained if we replace the empirical probability of error $\hat{L}_n(g_f)$ by a slightly larger quantity. To this end, let $\gamma > 0$ be a fixed parameter, and define the *margin error* by

$$L_n^\gamma(g_f) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{[f(X_i)(1-2Y_i) < \gamma]}.$$

Notice that for all $\gamma > 0$, $L_n^\gamma(g_f) \geq \hat{L}_n(g_f)$ and the $L_n^\gamma(g_f)$ is increasing in γ . An interpretation of the margin error $L_n^\gamma(g_f)$ is that it counts, apart from the number of misclassified pairs (X_i, Y_i) , also those which are well classified but only with a small “confidence” (or “margin”) by g_f .

The purpose of this section is to present a result of Freund, Schapire, Bartlett, and Lee (1998) which states that the margin error is always a good approximate upper bound for the probability of error, at least if γ is not too small. The elegant proof shown here is due to Koltchinskii and Panchenko (2002).

Theorem 1.15. *For every $\epsilon > 0$,*

$$\mathbb{P} \left\{ \sup_{f \in \mathcal{F}} (L(g_f) - L_n^\gamma(g_f)) > \frac{2\sqrt{2}}{\gamma} \sqrt{\frac{V_{\mathcal{C}} \log(n+1)}{n}} + \epsilon \right\} \leq e^{-2n\epsilon^2}.$$

Thus, with very high probability, the probability of error of any classifier g_f , $f \in \mathcal{F}$, may be simultaneously upper bounded by the sum

$$L_n^\gamma(g_f) + \frac{2\sqrt{2}}{\gamma} \sqrt{\frac{V_{\mathcal{C}} \log(n+1)}{n}}$$

plus a term of the order $n^{-1/2}$. Notice that, as γ grows, the first term of the sum increases, while the second decreases. The bound can be very useful whenever a classifier has a small margin error for a relatively large γ (i.e., if the classifier classifies the training data well with high “confidence”) since the second term only depends on the VC dimension of the small base class \mathcal{C} . As shown in the next section, the second term in the above sum may be replaced by $(c/\gamma)\sqrt{V_{\mathcal{C}}/n}$ for some universal constant c .

The proof of the theorem crucially uses the following simple lemma, called the “contraction principle”. Here we cite a version tailored for our needs. For the proof, see Ledoux and Talagrand (1991), pages 112–113.

Lemma 1.8. *Let $Z_1(f), \dots, Z_n(f)$ be arbitrary real-valued bounded random variables indexed by an abstract parameter f and let $\sigma_1, \dots, \sigma_n$ be independent symmetric sign variables, independent of the $Z_i(f)$'s (i.e., $\mathbb{P}\{\sigma_i = -1\} = \mathbb{P}\{\sigma_i = 1\} = 1/2$). If $\phi : \mathcal{R} \rightarrow \mathcal{R}$ is a Lipschitz function such that $|\phi(x) - \phi(y)| \leq |x - y|$ with $\phi(0) = 0$, then*

$$\mathbb{E} \sup_f \sum_{i=1}^n \sigma_i \phi(Z_i(f)) \leq \mathbb{E} \sup_f \sum_{i=1}^n \sigma_i Z_i(f).$$

PROOF OF THEOREM 1.15. For any $\gamma > 0$, introduce the function

$$\phi_\gamma(x) = \begin{cases} 1 & \text{if } x \leq 0 \\ 0 & \text{if } x \geq \gamma \\ 1 - x/\gamma & \text{if } x \in (0, \gamma) \end{cases}$$

Observe that $\mathbb{I}_{[x \leq 0]} \leq \phi_\gamma(x) \leq \mathbb{I}_{[x \leq \gamma]}$. Thus,

$$\sup_{f \in \mathcal{F}} (L(g_f) - L_n^\gamma(g_f)) \leq \sup_{f \in \mathcal{F}} \left(\mathbb{E} \phi_\gamma((1 - 2Y)f(X)) - \frac{1}{n} \sum_{i=1}^n \phi_\gamma((1 - 2Y_i)f(X)) \right).$$

Introduce the notation $Z(f) = (1 - 2Y)f(X)$ and $Z_i(f) = (1 - 2Y_i)f(X_i)$. Clearly, by the bounded difference inequality,

$$\begin{aligned} & \mathbb{P} \left\{ \sup_{f \in \mathcal{F}} \left(\mathbb{E} \phi_\gamma(Z(f)) - \frac{1}{n} \sum_{i=1}^n \phi_\gamma(Z_i(f)) \right) \right. \\ & \quad \left. > \mathbb{E} \sup_{f \in \mathcal{F}} \left(\mathbb{E} \phi_\gamma(Z(f)) - \frac{1}{n} \sum_{i=1}^n \phi_\gamma(Z_i(f)) \right) + \epsilon \right\} \leq e^{-2n\epsilon^2} \end{aligned}$$

and therefore it suffices to prove that the expected value of the supremum is bounded by $\frac{2\sqrt{2}}{\gamma} \sqrt{\frac{V_C \log(n+1)}{n}}$. As a first step, we proceed by a symmetrization argument just like in the proof of Theorem 1.9 to obtain

$$\begin{aligned} \mathbb{E} \sup_{f \in \mathcal{F}} \left(\mathbb{E} \phi_\gamma(Z(f)) - \frac{1}{n} \sum_{i=1}^n \phi_\gamma(Z_i(f)) \right) & \leq \mathbb{E} \sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{i=1}^n \sigma_i (\phi_\gamma(Z'_i(f)) - \phi_\gamma(Z_i(f))) \right) \\ & \leq 2\mathbb{E} \sup_{f \in \mathcal{F}} \left(\frac{1}{n} \sum_{i=1}^n \sigma_i (\phi_\gamma(Z_i(f)) - \phi_\gamma(0)) \right) \end{aligned}$$

where $\sigma_1, \dots, \sigma_n$ are i.i.d. symmetric sign variables and $Z'_i(f) = (1 - 2Y'_i)f(X'_i)$ where the (X'_i, Y'_i) are independent of the (X_i, Y_i) and have the same distribution as that of the pairs (X_i, Y_i) .

Observe that the function $\phi(x) = \gamma(\phi_\gamma(x) - \phi_\gamma(0))$ is Lipschitz and $\phi(0) = 0$, therefore, by the contraction principle (Lemma 1.8),

$$\mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (\phi_\gamma(Z_i(f)) - \phi_\gamma(0)) \leq \frac{1}{\gamma} \mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i Z_i(f) = \frac{1}{\gamma} \mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(X_i)$$

where at the last step we used the fact that $\sigma_i(1 - 2Y_i)$ is a symmetric sign variable, independent of the X_i and therefore $\sigma_i(1 - 2Y_i)f(X_i)$ has the same distribution as that of $\sigma_i f(X_i)$. The last expectation may be rewritten as

$$\mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(X_i) = \frac{1}{n} \mathbb{E} \sup_{N \geq 1} \sup_{g_1, \dots, g_N \in \mathcal{C}} \sup_{w_1, \dots, w_N} \sum_{i=1}^n \sum_{j=1}^N w_j \sigma_i g_j(X_i).$$

The key observation is that for any N and base classifiers g_i, \dots, g_N , the supremum in

$$\sup_{w_1, \dots, w_N} \sum_{i=1}^n \sum_{j=1}^N w_j \sigma_i g_j(X_i)$$

is achieved for a weight vector which puts all the mass in one index, that is, when $w_j = 1$ for some j . (This may be seen by observing that a linear function over a convex polygon achieves its maximum at one of the vertices of the polygon.) Thus,

$$\mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(X_i) = \frac{1}{n} \mathbb{E} \sup_{g \in \mathcal{C}} \sum_{i=1}^n \sigma_i g(X_i).$$

However, repeating the argument in the proof of Theorem 1.9 with the necessary adjustments, we obtain

$$\frac{1}{n} \mathbb{E} \sup_{g \in \mathcal{C}} \left| \sum_{i=1}^n \sigma_i g(X_i) \right| \leq \sqrt{\frac{2 \log \mathbb{S}_{\mathcal{C}}(n)}{n}} \leq \sqrt{\frac{2V_{\mathcal{C}} \log(n+1)}{n}}$$

which completes the proof of the desired inequality. \square

1.4.6 Appendix: sharper bounds via chaining

In this section we present an improvement of the Vapnik-Chervonenkis inequality stating that for any class \mathcal{A} of sets of VC dimension V ,

$$\mathbb{E} \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \leq c \sqrt{\frac{V}{n}},$$

where c is a universal constant. This in turn implies for empirical risk minimization that

$$\mathbb{E} L(g_n^*) - L_{\mathcal{C}} \leq 2c \sqrt{\frac{V_{\mathcal{C}}}{n}}.$$

The new bound involves some geometric and combinatorial quantities related to the class \mathcal{A} . Consider a pair of bit vectors $b = (b_1, \dots, b_n)$ and $c = (c_1, \dots, c_n)$ from $\{0, 1\}^n$, and define their distance by

$$\rho(b, c) = \sqrt{\frac{1}{n} \sum_{i=1}^n \mathbb{I}_{[b_i \neq c_i]}}.$$

Thus, $\rho(b, c)$ is just the square root of the normalized Hamming distance between b and c . Observe that ρ may also be considered as the normalized euclidean distance between the corners of the hypercube $[0, 1]^n \subset \mathcal{R}^n$, and therefore it is indeed a distance.

Now let $B \subset \{0, 1\}^n$ be any set of bit vectors, and define a *cover* of radius $r > 0$ as a set $B_r \subset \{0, 1\}^n$ such that for any $b \in B$ there exists a $c \in B_r$ such that $\rho(b, c) \leq r$. The *covering number* $N(r, B)$ is the cardinality of the smallest cover of radius r .

A class \mathcal{A} of subsets of \mathcal{R}^d and a set of n points $x_1^n = \{x_1, \dots, x_n\} \subset \mathcal{R}^d$ define a set of bit vectors by

$$\mathcal{A}(x_1^n) = \{b = (b_1, \dots, b_n) \in \{0, 1\}^n : b_i = \mathbb{I}_{[x_i \in A]}, i = 1, \dots, n \text{ for some } A \in \mathcal{A}\}.$$

That is, every bit vector $b \in \mathcal{A}(x_1^n)$ describes the intersection of $\{x_1, \dots, x_n\}$ with a set A in \mathcal{A} . We have the following:

Theorem 1.16.

$$\mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right\} \leq \frac{24}{\sqrt{n}} \max_{x_1, \dots, x_n \in \mathcal{R}^d} \int_0^1 \sqrt{\log 2N(r, \mathcal{A}(x_1^n))} dr.$$

The theorem implies that $\mathbb{E} \{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \} = O(1/\sqrt{n})$ whenever the integral in the bound is uniformly bounded over all x_1, \dots, x_n and all n . Note that the bound of Theorem 1.9 is always of larger order of magnitude, trivial cases excepted. The main additional idea is Dudley's *chaining* trick.

PROOF. As in the proof of Theorem 1.9, we see that

$$\begin{aligned}
& \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right\} \\
& \leq \frac{1}{n} \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbb{I}_{[X_i \in A]} - \mathbb{I}_{[X'_i \in A]}) \right| \right\} \\
& \leq \frac{1}{n} \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i \mathbb{I}_{[X_i \in A]} \right| \right\} + \frac{1}{n} \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i \mathbb{I}_{[X'_i \in A]} \right| \right\} \\
& = \frac{2}{n} \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i \mathbb{I}_{[X_i \in A]} \right| \right\} \\
& = \frac{2}{n} \mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i \mathbb{I}_{[X_i \in A]} \right| \middle| X_1, \dots, X_n \right\}.
\end{aligned}$$

Just as in the proof of theorem 1, we fix the values $X_1 = x_1, \dots, X_n = x_n$ and study

$$\mathbb{E} \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i \mathbb{I}_{[x_i \in A]} \right| \right\} = \mathbb{E} \left\{ \max_{b \in \mathcal{A}(x_1^n)} \left| \sum_{i=1}^n \sigma_i b_i \right| \right\}.$$

Now let $B_0 \stackrel{\text{def}}{=} \{b^{(0)}\}$ be the singleton set containing the all-zero vector $b^{(0)} = (0, \dots, 0)$, and let B_1, B_2, \dots, B_M be subsets of $\{0, 1\}^n$ such that each B_k is a minimal cover of $\mathcal{A}(x_1^n)$ of radius 2^{-k} , and $M = \lfloor \log_2 \sqrt{n} \rfloor + 1$. Note that B_0 is also a cover of radius 2^0 , and that $B_M = \mathcal{A}(x_1^n)$. Now denote the (random) vector reaching the maximum by $b^* = (b_1^*, \dots, b_n^*) \in \mathcal{A}(x_1^n)$, that is,

$$\left| \sum_{i=1}^n \sigma_i b_i^* \right| = \max_{b \in \mathcal{A}(x_1^n)} \left| \sum_{i=1}^n \sigma_i b_i \right|,$$

and, for each $k \leq M$, let $b^{(k)} \in B_k$ be a nearest neighbor of b^* in the k -th cover, that is,

$$\rho(b^{(k)}, b^*) \leq \rho(b, b^*) \quad \text{for all } b \in B_k.$$

Note that $\rho(b^{(k)}, b^*) \leq 2^{-k}$, and therefore

$$\rho(b^{(k)}, b^{(k-1)}) \leq \rho(b^{(k)}, b^*) + \rho(b^{(k-1)}, b^*) \leq 3 \cdot 2^{-k}.$$

Now clearly,

$$\begin{aligned}
\sum_{i=1}^n \sigma_i b_i^* &= \sum_{i=1}^n \sigma_i b_i^{(0)} + \sum_{k=1}^M \sum_{i=1}^n \sigma_i (b_i^{(k)} - b_i^{(k-1)}) \\
&= \sum_{k=1}^M \sum_{i=1}^n \sigma_i (b_i^{(k)} - b_i^{(k-1)}),
\end{aligned}$$

so

$$\begin{aligned} \mathbb{E} \left\{ \max_{b \in \mathcal{A}(x_1^n)} \left| \sum_{i=1}^n \sigma_i b_i \right| \right\} &= \mathbb{E} \left| \sum_{k=1}^M \sum_{i=1}^n \sigma_i (b_i^{(k)} - b_i^{(k-1)}) \right| \\ &\leq \sum_{k=1}^M \mathbb{E} \left| \sum_{i=1}^n \sigma_i (b_i^{(k)} - b_i^{(k-1)}) \right| \\ &\leq \sum_{k=1}^M \mathbb{E} \max_{b \in B_k, c \in B_{k-1}: \rho(b,c) \leq 3 \cdot 2^{-k}} \left| \sum_{i=1}^n \sigma_i (b_i - c_i) \right|. \end{aligned}$$

Now it follows from Lemma 1.2 that for each pair $b \in B_k, c \in B_{k-1}$ with $\rho(b, c) \leq 3 \cdot 2^{-k}$, and for all $s > 0$,

$$e^{s \sum_{i=1}^n \sigma_i (b_i - c_i)} \leq e^{s^2 n (3 \cdot 2^{-k})^2 / 2}.$$

On the other hand, the number of such pairs is bounded by $|B_k| \cdot |B_{k-1}| \leq |B_k|^2 = N(2^{-k}, \mathcal{A}(x_1^n))^2$. Then Lemma 1.3 implies that for each $1 \leq k \leq M$,

$$\mathbb{E} \max_{b \in B_k, c \in B_{k-1}: \rho(b,c) \leq 3 \cdot 2^{-k}} \left| \sum_{i=1}^n \sigma_i (b_i - c_i) \right| \leq 3\sqrt{n} 2^{-k} \sqrt{2 \log 2N(2^{-k}, \mathcal{A}(x_1^n))^2}.$$

Summarizing, we obtain

$$\begin{aligned} \mathbb{E} \left\{ \max_{b \in \mathcal{A}(x_1^n)} \left| \sum_{i=1}^n \sigma_i b_i \right| \right\} &\leq 3\sqrt{n} \sum_{k=1}^M 2^{-k} \sqrt{2 \log 2N(2^{-k}, \mathcal{A}(x_1^n))^2} \\ &\leq 12\sqrt{n} \sum_{k=1}^{\infty} 2^{-(k+1)} \sqrt{\log 2N(2^{-k}, \mathcal{A}(x_1^n))} \\ &\leq 12\sqrt{n} \int_0^1 \sqrt{\log 2N(r, \mathcal{A}(x_1^n))} dr, \end{aligned}$$

where at the last step we used the fact that $N(r, \mathcal{A}(x_1^n))$ is a monotonically decreasing function of r . The proof is finished. \square

To complete our argument, we need to relate the vc dimension of a class of sets \mathcal{A} to the covering numbers $N(r, \mathcal{A}(x_1^n))$ appearing in Theorem 3.10.

Theorem 1.17. *Let \mathcal{A} be a class of sets with vc dimension $V < \infty$. For every $x_1, \dots, x_n \in \mathcal{R}^d$ and $0 \leq r \leq 1$,*

$$N(r, \mathcal{A}(x_1^n)) \leq \left(\frac{4e}{r^2} \right)^{V/(1-1/e)}.$$

Theorem 1.17 is due to Dudley (1978). Haussler (1995) refined Dudley's probabilistic argument and showed that the stronger bound

$$N(r, \mathcal{A}(x_1^n)) \leq e(V+1) \left(\frac{2e}{r^2} \right)^V.$$

also holds.

PROOF. Fix x_1, \dots, x_n , and consider the set $B_0 = \mathcal{A}(x_1^n) \in \{0, 1\}^n$. Fix $r \in (0, 1)$, and let $B_r \subset \{0, 1\}^n$ be a minimal cover of B_0 of radius r with respect to the metric

$$\rho(b, c) = \sqrt{\frac{1}{n} \sum_{i=1}^n \mathbb{I}_{[b_i \neq c_i]}}.$$

We need to show that $|B_r| \leq (4e/r^2)^{V/(1-1/e)}$.

First note that there exists a “packing set” $C_r \subset B_0$ such that $|B_r| \leq |C_r|$ and any two elements $b, c \in C_r$ are r -separated, that is, $\rho(b, c) > r$. To see this, suppose that C_r is such an r -separated set of maximal cardinality. Then for any $b \in B_0$, there exists a $c \in C_r$ with $\rho(b, c) \leq r$, since otherwise adding b to the set C_r would increase its cardinality, and it would still be r -separated. Thus, C_r is a cover of radius r , which implies that $|B_r| \leq |C_r|$. Denote the elements of C_r by $c^{(1)}, \dots, c^{(M)}$, where $M = |C_r|$. For any $i, j \leq M$, define $A_{i,j}$ as the set of indices where the binary vectors $c^{(i)}$ and $c^{(j)}$ disagree:

$$A_{i,j} = \left\{ 1 \leq m \leq n : c_m^{(i)} \neq c_m^{(j)} \right\}.$$

Note that any two elements of C_r differ in at least nr^2 components. Next define K independent random variables Y_1, \dots, Y_K , distributed uniformly over the set $\{1, 2, \dots, n\}$, where K will be specified later. Then for any $i, j \leq M$, $i \neq j$, and $k \leq K$,

$$\mathbb{P}\{Y_k \in A_{i,j}\} \geq r^2,$$

and therefore the probability that no one of Y_1, \dots, Y_K falls in the set $A_{i,j}$ is less than $(1 - r^2)^K$. Observing that there are less than M^2 sets $A_{i,j}$, and applying the union bound, we obtain that

$$\begin{aligned} &\mathbb{P}\{\text{for all } i \neq j, i, j \leq M, \text{ at least one } Y_k \text{ falls in } A_{i,j}\} \\ &\geq 1 - M^2(1 - r^2)^K \geq 1 - M^2 e^{-Kr^2}. \end{aligned}$$

If we choose $K = \lceil 2 \log M/r^2 \rceil + 1$, then the above probability is strictly positive. This implies that there exist $K = \lceil 2 \log M/r^2 \rceil + 1$ indices $y_1, \dots, y_K \in \{1, 2, \dots, n\}$ such that at least one y_k falls in each set $A_{i,j}$. Therefore, restricted to the K components y_1, \dots, y_K ,

the elements of C_r are all different, and since $C_r \subset B_0$, C_r does not shatter any set of size larger than V . Therefore, by Sauer's lemma we obtain

$$|C_r| = M \leq \left(\frac{eK}{V} \right)^V$$

for $K \leq V$. Thus, if $\log M \geq V$, then

$$\begin{aligned} \log M &\leq V \log \frac{e(\lceil 2\log M/r^2 \rceil + 1)}{V} \\ &\leq V \left(\log \frac{4e}{r^2} + \log \frac{\log M}{V} \right) \\ &\leq V \log \frac{4e}{r^2} + \frac{1}{e} \log M \quad (\text{since } \log x \leq x/e \text{ for } x > 0). \end{aligned}$$

Therefore,

$$\log M \leq \frac{V}{1 - 1/e} \log \frac{4e}{r^2}.$$

If $\log M < V$, then the above inequality holds trivially. This concludes the proof. \square

Combining this result with Theorem 3.10 we obtain that for any class \mathcal{A} with vc dimension V ,

$$\mathbb{E} \left\{ \sup_{A \in \mathcal{A}} |\mu_n(A) - \mu(A)| \right\} \leq c \sqrt{\frac{V}{n}},$$

where c is a universal constant.

1.5 Minimax lower bounds

The purpose of this section is to investigate how good the bounds obtained in the previous chapter for empirical risk minimization are. We have seen that for any class \mathcal{C} of classifiers with vc dimension V , a classifier g_n^* minimizing the empirical risk satisfies

$$\mathbb{E}L(g_n^*) - L_c \leq O \left(\sqrt{\frac{L_c V_c \log n}{n}} + \frac{V_c \log n}{n} \right),$$

and also

$$\mathbb{E}L(g_n^*) - L_c \leq O \left(\sqrt{\frac{V_c}{n}} \right).$$

In this section we seek answers for the following questions: Are these upper bounds (at least up to the order of magnitude) tight? Is there a much better way of selecting a classifier than minimizing the empirical error?

Let us formulate exactly what we are interested in. Let \mathcal{C} be a class of decision functions $g : \mathcal{R}^d \rightarrow \{0, 1\}$. The training sequence $D_n = ((X_1, Y_1), \dots, (X_n, Y_n))$ is used to select the classifier $g_n(X) = g_n(X, D_n)$ from \mathcal{C} , where the selection is based on the data D_n . We emphasize here that g_n can be an arbitrary function of the data, we do not restrict our attention to empirical error minimization, where g_n is a classifier in \mathcal{C} that minimizes the number errors committed on the data D_n .

As before, we measure the performance of the selected classifier by the difference between the error probability $L(g_n) = \mathbb{P}\{g_n(X) \neq Y | D_n\}$ of the selected classifier and that of the best in the class, $L_{\mathcal{C}}$. In particular, we seek lower bounds for

$$\sup \mathbb{E}L(g_n) - L_{\mathcal{C}},$$

where the supremum is taken over all possible distributions of the pair (X, Y) . A lower bound for this quantities means that no matter what our method of picking a rule from \mathcal{C} is, we may face a distribution such that our method performs worse than the bound.

Actually, we investigate a stronger problem, in that the supremum is taken over all distributions with $L_{\mathcal{C}}$ kept at a fixed value between zero and $1/2$. We will see that the bounds depend on n , $V_{\mathcal{C}}$, and $L_{\mathcal{C}}$ jointly. As it turns out, the situations for $L_{\mathcal{C}} > 0$ and $L_{\mathcal{C}} = 0$ are quite different. Because of its simplicity, we first treat the case $L_{\mathcal{C}} = 0$. All the proofs are based on a technique called “the probabilistic method.” The basic idea here is that the existence of a “bad” distribution is proved by considering a large class of distributions, and bounding the average behavior over the class.

1.5.1 The zero-error case

Here we obtain lower bounds under the assumption that the best classifier in the class has zero error probability. Recall that by Corollary 1.2 the expected probability of error of an empirical risk minimizer is bounded by $O(V_{\mathcal{C}} \log n / n)$. Next we obtain minimax lower bounds close to the upper bounds.

Theorem 1.18. *Let \mathcal{C} be a class of discrimination functions with VC dimension V . Let \mathcal{X} be the set of all random variables (X, Y) for which $L_{\mathcal{C}} = 0$. Then, for every discrimination rule g_n based upon $X_1, Y_1, \dots, X_n, Y_n$, and $n \geq V - 1$,*

$$\sup_{(X, Y) \in \mathcal{X}} \mathbb{E}L(g_n) \geq \frac{V - 1}{2en} \left(1 - \frac{1}{n}\right).$$

PROOF. The idea is to construct a family \mathcal{F} of 2^{V-1} distributions within the distributions with $L_{\mathcal{C}} = 0$ as follows: first find points x_1, \dots, x_V that are shattered by \mathcal{C} . Each distribution in \mathcal{F} is concentrated on the set of these points. A member in \mathcal{F} is described by $V - 1$ bits,

b_1, \dots, b_{V-1} . For convenience, this is represented as a bit vector b . Assume $V - 1 \leq n$. For a particular bit vector, we let $X = x_i$ ($i < V$) with probability $1/n$ each, while $X = x_V$ with probability $1 - (V - 1)/n$. Then set $Y = f_b(X)$, where f_b is defined as follows:

$$f_b(x) = \begin{cases} b_i & \text{if } x = x_i, i < V \\ 0 & \text{if } x = x_V. \end{cases}$$

Note that since Y is a function of X , we must have $L^* = 0$. Also, $L_C = 0$, as the set $\{x_1, \dots, x_V\}$ is shattered by \mathcal{C} , i.e., there is a $g \in \mathcal{C}$ with $g(x_i) = f_b(x_i)$ for $1 \leq i \leq V$. Clearly,

$$\begin{aligned} & \sup_{(X,Y):L_C=0} \mathbb{E}\{L(g_n) - L_C\} \\ & \geq \sup_{(X,Y) \in \mathcal{F}} \mathbb{E}\{L(g_n) - L_C\} \\ & = \sup_b \mathbb{E}\{L(g_n) - L_C\} \\ & \geq \mathbb{E}\{L(g_n) - L_C\} \\ & \quad (\text{where } b \text{ is replaced by } B, \text{ uniformly distributed over } \{0,1\}^{V-1}) \\ & = \mathbb{E}\{L(g_n)\}, \\ & = \mathbb{P}\{g_n(X, X_1, Y_1, \dots, X_n, Y_n) \neq f_B(X)\}. \end{aligned}$$

The last probability may be viewed as the error probability of the decision function $g_n : \mathcal{R}^d \times (\mathcal{R}^d \times \{0,1\})^n \rightarrow \{0,1\}$ in predicting the value of the random variable $f_B(X)$ based on the observation $Z_n = (X, X_1, Y_1, \dots, X_n, Y_n)$. Naturally, this probability is bounded from below by the Bayes probability of error

$$L^*(Z_n, f_B(X)) = \inf_{g_n} \mathbb{P}\{g_n(Z_n) \neq f_B(X)\}$$

corresponding to the decision problem $(Z_n, f_B(X))$. By the results of Chapter 1,

$$L^*(Z_n, f_B(X)) = \mathbb{E}\{\min(\eta^*(Z_n), 1 - \eta^*(Z_n))\},$$

where $\eta^*(Z_n) = \mathbb{P}\{f_B(X) = 1 | Z_n\}$. Observe that

$$\eta^*(Z_n) = \begin{cases} 1/2 & \text{if } X \neq X_1, \dots, X \neq X_n, X \neq x_V \\ 0 \text{ or } 1 & \text{otherwise.} \end{cases}$$

Thus, we see that

$$\begin{aligned}
\sup_{(X,Y):L_C=0} \mathbb{E}\{L(g_n) - L_C\} &\geq L^*(Z_n, f_B(X)) \\
&= \frac{1}{2} \mathbb{P}\{X \neq X_1, \dots, X \neq X_n, X \neq x_V\} \\
&= \frac{1}{2} \sum_{i=1}^{V-1} \mathbb{P}\{X = x_i\} (1 - \mathbb{P}\{X = x_i\})^n \\
&= \frac{V-1}{2n} (1 - 1/n)^n \\
&\geq \frac{V-1}{2en} \left(1 - \frac{1}{n}\right) \quad (\text{since } (1 - 1/n)^{n-1} \downarrow 1/e).
\end{aligned}$$

This concludes the proof. \square

1.5.2 The general case

In the more general case, when the best decision in the class \mathcal{C} has positive error probability, the upper bounds derived in Chapter 2 for the expected error probability of the classifier obtained by minimizing the empirical risk are much larger than when $L_C = 0$. Theorem 1.19 below gives a lower bound for $\sup_{(X,Y):L_C \text{ fixed}} \mathbb{E}L(g_n) - L_C$. As a function of n and V_C , the bound decreases basically as in the upper bound obtained from Theorem 1.11. Interestingly, the lower bound becomes smaller as L_C decreases, as should be expected. The bound is largest when L_C is close to 1/2.

Theorem 1.19. *Let \mathcal{C} be a class of discrimination functions with VC dimension $V \geq 2$. Let \mathcal{X} be the set of all random variables (X, Y) for which for fixed $L \in (0, 1/2)$,*

$$L = \inf_{g \in \mathcal{C}} \mathbb{P}\{g(X) \neq Y\}.$$

Then, for every discrimination rule g_n based upon $X_1, Y_1, \dots, X_n, Y_n$,

$$\sup_{(X,Y) \in \mathcal{X}} \mathbb{E}(L(g_n) - L) \geq \sqrt{\frac{L(V-1)}{24n}} e^{-8} \quad \text{if } n \geq \frac{V-1}{2L} \max(9, 1/(1-2L)^2).$$

PROOF. Again we consider the finite family \mathcal{F} from the previous section. The notation b and B is also as above. X now puts mass p at x_i , $i < V$, and mass $1 - (V-1)p$ at x_V . This imposes the condition $(V-1)p \leq 1$, which will be satisfied. Next introduce the constant $c \in (0, 1/2)$. We no longer have Y as a function of X . Instead, we have a uniform $[0, 1]$

random variable U independent of X and define

$$Y = \begin{cases} 1 & \text{if } U \leq \frac{1}{2} - c + 2cb_i, X = x_i, i < V \\ 0 & \text{otherwise.} \end{cases}$$

Thus, when $X = x_i, i < V$, Y is 1 with probability $1/2 - c$ or $1/2 + c$. A simple argument shows that the best rule for b is the one which sets

$$f_b(x) = \begin{cases} 1 & \text{if } x = x_i, i < V, b_i = 1 \\ 0 & \text{otherwise.} \end{cases}$$

Also, observe that

$$L = (V - 1)p(1/2 - c).$$

Noting that $|2\eta(x_i) - 1| = c$ for $i < V$, for fixed b , we may write

$$L(g_n) - L \geq \sum_{i=1}^{V-1} 2pcI_{\{g_n(x_i, X_1, Y_1, \dots, X_n, Y_n) = 1 - f_b(x_i)\}}.$$

It is sometimes convenient to make the dependence of g_n upon b explicit by considering $g_n(x_i)$ as a function of $x_i, X_1, \dots, X_n, U_1, \dots, U_n$ (an i.i.d. sequence of uniform $[0, 1]$ random variables), and b_i . We replace b by a uniformly distributed random B over $\{0, 1\}^{V-1}$. After this randomization, denote $Z_n = (X, X_1, Y_1, \dots, X_n, Y_n)$. Thus,

$$\begin{aligned} \sup_{(X,Y) \in \mathcal{F}} \mathbb{E}\{L(g_n) - L\} &= \sup_b \mathbb{E}\{L(g_n) - L\} \\ &\geq \mathbb{E}\{L(g_n) - L\} \quad (\text{with random } B) \\ &\geq \sum_{i=1}^{V-1} 2pc\mathbb{E}I_{\{g_n(x_i, X_1, \dots, Y_n) = 1 - f_B(x_i)\}} \\ &= 2c\mathbb{P}\{g_n(Z_n) \neq f_B(X)\} \\ &\geq 2cL^*(Z_n, f_B(X)), \end{aligned}$$

where, as before, $L^*(Z_n, f_B(X))$ denotes the Bayes probability of error of predicting the value of $f_B(X)$ based on observing Z_n . All we have to do is to find a suitable lower bound for

$$L^*(Z_n, f_B(X)) = \mathbb{E}\{\min(\eta^*(Z_n), 1 - \eta^*(Z_n))\},$$

where $\eta^*(Z_n) = \mathbb{P}\{f_B(X) = 1 | Z_n\}$. Observe that

$$\eta^*(Z_n) = \begin{cases} 1/2 & \text{if } X \neq X_1, \dots, X \neq X_n \text{ and } X \neq x_V \\ \mathbb{P}\{B_i = 1 | Y_{i_1}, \dots, Y_{i_k}\} & \text{if } X = X_{i_1} = \dots = X_{i_k} = x_i, i < V. \end{cases}$$

Next we compute $\mathbb{P}\{B_i = 1|Y_{i_1} = y_1, \dots, Y_{i_k} = y_k\}$ for $y_1, \dots, y_k \in \{0, 1\}$. Denoting the numbers of zeros and ones by $k_0 = |\{j \leq k : y_j = 0\}|$ and $k_1 = |\{j \leq k : y_j = 1\}|$, we see that

$$\begin{aligned} & \mathbb{P}\{B_i = 1|Y_{i_1} = y_1, \dots, Y_{i_k} = y_k\} \\ &= \frac{(1 - 2c)^{k_1}(1 + 2c)^{k_0}}{(1 - 2c)^{k_1}(1 + 2c)^{k_0} + (1 + 2c)^{k_1}(1 - 2c)^{k_0}}. \end{aligned}$$

Therefore, if $X = X_{i_1} = \dots = X_{i_k} = x_i$, $i < V$, then

$$\begin{aligned} & \min(\eta^*(Z_n), 1 - \eta^*(Z_n)) \\ &= \frac{\min((1 - 2c)^{k_1}(1 + 2c)^{k_0}, (1 + 2c)^{k_1}(1 - 2c)^{k_0})}{(1 - 2c)^{k_1}(1 + 2c)^{k_0} + (1 + 2c)^{k_1}(1 - 2c)^{k_0}} \\ &= \frac{\min\left(1, \left(\frac{1+2c}{1-2c}\right)^{k_1-k_0}\right)}{1 + \left(\frac{1+2c}{1-2c}\right)^{k_1-k_0}} \\ &= \frac{1}{1 + \left(\frac{1+2c}{1-2c}\right)^{|k_1-k_0}|}. \end{aligned}$$

In summary, denoting $a = (1 + 2c)/(1 - 2c)$, we have

$$\begin{aligned} L^*(Z_n, f_B(X)) &= \mathbb{E}\left\{\frac{1}{1+a^{\sum_{j:X_j=x}(2Y_j-1)}}\right\} \\ &\geq \mathbb{E}\left\{\frac{1}{2a^{\sum_{j:X_j=x}(2Y_j-1)}}\right\} \\ &\geq \frac{1}{2} \sum_{i=1}^{V-1} \mathbb{P}\{X = x_i\} \mathbb{E}\left\{a^{-\sum_{j:X_j=x_i}(2Y_j-1)}\right\} \\ &\geq \frac{1}{2}(V-1)pa^{-\mathbb{E}\{\sum_{j:X_j=x_i}(2Y_j-1)\}} \\ &\quad (\text{by Jensen's inequality}). \end{aligned}$$

Next we bound $\mathbb{E}\left\{\left|\sum_{j:X_j=x_i}(2Y_j-1)\right|\right\}$. Clearly, if $B(k, q)$ denotes a binomial random variable with parameters k and q ,

$$\mathbb{E}\left\{\left|\sum_{j:X_j=x_i}(2Y_j-1)\right|\right\} = \sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} \mathbb{E}\{|2B(k, 1/2 - c) - k|\}.$$

However, by straightforward calculation we see that

$$\begin{aligned}\mathbb{E}\{|2B(k, 1/2 - c) - k|\} &\leq \sqrt{\mathbb{E}\{(2B(k, 1/2 - c) - k)^2\}} \\ &= \sqrt{k(1 - 4c^2) + 4k^2c^2} \\ &\leq 2kc + \sqrt{k}.\end{aligned}$$

Therefore, applying Jensen's inequality once again, we get

$$\sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} \mathbb{E}\{|2B(k, 1/2 - c) - k|\} \leq 2npc + \sqrt{np}.$$

Summarizing what we have obtained so far, we have

$$\begin{aligned}\sup_b \mathbb{E}\{L(g_n) - L\} &\geq 2cL^*(Z_n, f_B(X)) \\ &\geq 2c \frac{1}{2}(V - 1)p e^{-2npc - \sqrt{np}} \\ &\geq c(V - 1)p e^{-2npc(a-1) - (a-1)\sqrt{np}} \\ &\quad (\text{by the inequality } 1 + x \leq e^x) \\ &= c(V - 1)p e^{-8npc^2/(1-2c) - 4c\sqrt{np}/(1-2c)}.\end{aligned}$$

A rough asymptotic analysis shows that the best asymptotic choice for c is given by

$$c = \frac{1}{\sqrt{4np}}.$$

Then the constraint $L = (V - 1)p(1/2 - c)$ leaves us with a quadratic equation in c . Instead of solving this equation, it is more convenient to take $c = \sqrt{(V - 1)/(8nL)}$. If $2nL/(V - 1) \geq 9$, then $c \leq 1/6$. With this choice for c , using $L = (V - 1)p(1/2 - c)$, straightforward calculation provides

$$\sup_{(X,Y) \in \mathcal{F}} \mathbb{E}(L(g_n) - L) \geq \sqrt{\frac{(V - 1)L}{24n}} e^{-8}.$$

The condition $p(V - 1) \leq 1$ implies that we need to ask that $n \geq (V - 1)/(2L(1 - 2L)^2)$. This concludes the proof of Theorem 1.19. \square

1.6 Complexity regularization

This section deals with the problem of automatic model selection. Our goal is to develop some data-based methods to find the class \mathcal{C} of classifiers in a way that approximately minimizes the probability of error of the empirical risk minimizer.

1.6.1 Model selection by penalization

In empirical risk minimization one selects a classifier from a given class \mathcal{C} by minimizing the error estimate $\hat{L}_n(g)$ over all $g \in \mathcal{C}$. This provides an estimate whose loss is close to the optimal loss L^* if the class \mathcal{C} is (i) sufficiently large so that the loss of the best function in \mathcal{C} is close to L^* and (ii) is sufficiently small so that finding the best candidate in \mathcal{C} based on the data is still possible. These two requirements are clearly in conflict. The trade-off is best understood by writing

$$\mathbb{E}L(g_n^*) - L^* = \left(\mathbb{E}L(g_n^*) - \inf_{g \in \mathcal{C}} L(g) \right) + \left(\inf_{g \in \mathcal{C}} L(g) - L^* \right).$$

The first term is often called *estimation error*, while the second is the *approximation error*.

It is common to fix in advance a sequence of model classes $\mathcal{C}_1, \mathcal{C}_2, \dots$, which, typically, become richer for larger indices. Given the data D_n , one wishes to select a good model from *one* of these classes. This is the problem of model selection.

Denote by \hat{g}_k a function in \mathcal{C}_k having minimal empirical risk. One hopes to select a model class \mathcal{C}_K such that the excess error $\mathbb{E}L(\hat{g}_K) - L^*$ is close to

$$\min_k \mathbb{E}L(\hat{g}_k) - L^* = \min_k \left[\left(\mathbb{E}L(\hat{g}_k) - \inf_{g \in \mathcal{C}_k} L(g) \right) + \left(\inf_{g \in \mathcal{C}_k} L(g) - L^* \right) \right].$$

The idea of *structural risk minimization*, (also known as *complexity regularization*), is to add a complexity penalty to each of the $\hat{L}_n(\hat{g}_k)$'s to compensate for the overfitting effect. This penalty is usually closely related to a distribution-free upper bound for $\sup_{g \in \mathcal{C}_k} |\hat{L}_n(g) - L(g)|$ so that the penalty eliminates the effect of overfitting.

The first general result shows that any approximate upper bound on error can be used to define a (possibly data-dependent) complexity penalty $C_n(k)$ and a model selection algorithm for which the excess error is close to

$$\min_k \left[\mathbb{E}C_n(k) + \left(\inf_{g \in \mathcal{C}_k} L(g) - L^* \right) \right].$$

Our goal is to select, among the classifiers \hat{g}_k one which has approximately minimal loss. The key assumption for our analysis is that the true loss of \hat{g}_k can be estimated for all k .

Assumption 1 *There are positive numbers c and m such that for each k an estimate $R_{n,k}$ on $L(\hat{g}_k)$ is available which satisfies*

$$\mathbb{P}[L(\hat{g}_k) > R_{n,k} + \epsilon] \leq ce^{-2m\epsilon^2}$$

for all $\epsilon > 0$.

Now define the complexity penalty by

$$C_n(k) = R_{n,k} - \hat{L}_n(\hat{g}_k) + \sqrt{\frac{\log k}{m}}.$$

The last term is required because of technical reasons that will become apparent shortly. It is typically small. The difference $R_{n,k} - \hat{L}_n(\hat{g}_k)$ is simply an estimate of the ‘right’ amount of penalization $L(\hat{g}_k) - \hat{L}_n(\hat{g}_k)$. Finally, define the prediction rule:

$$g_n^* = \arg \min_{k=1,2,\dots} \tilde{L}_n(\hat{g}_k),$$

where

$$\tilde{L}_n(\hat{g}_k) = \hat{L}_n(\hat{g}_k) + C_n(k) = R_{n,k} + \sqrt{\frac{\log k}{m}}.$$

The following theorem summarizes the main performance bound for g_n^* .

Theorem 1.20. *Assume that the error estimates $R_{n,k}$ satisfy Assumption 1 for some positive constants c and m . Then*

$$\mathbb{E}L(g_n^*) - L^* \leq \min_k \left[\mathbb{E}C_n(k) + \left(\inf_{g \in \mathcal{C}_k} L(g) - L^* \right) \right] + \sqrt{\frac{\log(ce)}{2m}}.$$

Theorem 1.20 shows that the prediction rule minimizing the penalized empirical loss achieves an almost optimal trade-off between the approximation error and the expected complexity, provided that the estimate $R_{n,k}$ on which the complexity is based is an approximate upper bound on the loss. In particular, if we knew in advance which of the classes \mathcal{C}_k contained the optimal prediction rule, we could use the error estimates $R_{n,k}$ to obtain an upper bound on $\mathbb{E}L(\hat{g}_k) - L^*$, and this upper bound would not improve on the bound of Theorem 1.20 by more than $O\left(\sqrt{\log k/m}\right)$.

PROOF. For brevity, introduce the notation

$$L_k^* = \inf_{g \in \mathcal{C}_k} L(g).$$

Then for any $\epsilon > 0$,

$$\begin{aligned}
\mathbb{P} \left[L(g_n^*) - \tilde{L}_n(g_n^*) > \epsilon \right] &\leq \mathbb{P} \left[\sup_{j=1,2,\dots} \left(L(\hat{g}_j) - \tilde{L}_n(\hat{g}_j) \right) > \epsilon \right] \\
&\leq \sum_{j=1}^{\infty} \mathbb{P} \left[L(\hat{g}_j) - \tilde{L}_n(\hat{g}_j) > \epsilon \right] \\
&\quad (\text{by the union bound}) \\
&= \sum_{j=1}^{\infty} \mathbb{P} \left[L(\hat{g}_j) - R_{n,j} > \epsilon + \sqrt{\frac{\log j}{m}} \right] \\
&\quad (\text{by definition}) \\
&\leq \sum_{j=1}^{\infty} ce^{-2m(\epsilon + \sqrt{\frac{\log j}{m}})^2} \quad (\text{by Assumption 1}) \\
&\leq \sum_{j=1}^{\infty} ce^{-2m(\epsilon^2 + \frac{\log j}{m})} \\
&< 2ce^{-2m\epsilon^2} \quad (\text{since } \sum_{j=1}^{\infty} j^{-2} < 2).
\end{aligned}$$

To prove the theorem, for each k , we decompose $L(g_n^*) - L_k^*$ as

$$L(g_n^*) - L_k^* = \left(L(g_n^*) - \inf_j \tilde{L}_n(\hat{g}_j) \right) + \left(\inf_j \tilde{L}_n(\hat{g}_j) - L_k^* \right).$$

The first term may be bounded, by standard integration of the tail inequality shown above, as $\mathbb{E} \left[L(g_n^*) - \inf_j \tilde{L}_n(\hat{g}_j) \right] \leq \sqrt{\log(ce)/(2m)}$. Choosing g_k^* such that $L(g_k^*) = L_k^*$, the second term may be bounded directly by

$$\begin{aligned}
\mathbb{E} \inf_j \tilde{L}_n(\hat{g}_j) - L_k^* &\leq \mathbb{E} \tilde{L}_n(\hat{g}_k) - L_k^* \\
&= \mathbb{E} \hat{L}_n(\hat{g}_k) - L_k^* + \mathbb{E} C_n(k) \\
&\quad (\text{by the definition of } \tilde{L}_n(\hat{g}_k)) \\
&\leq \mathbb{E} \hat{L}_n(g_k^*) - L(g_k^*) + \mathbb{E} C_n(k) \\
&\quad (\text{since } \hat{g}_k \text{ minimizes the empirical loss on } \mathcal{C}_k) \\
&= \mathbb{E} C_n(k),
\end{aligned}$$

where the last step follows from the fact that $\mathbb{E} \hat{L}_n(g_k^*) = L(g_k^*)$. Summing the obtained bounds for both terms yields that for each k ,

$$\mathbb{E} L(g_n^*) \leq \mathbb{E} C_n(k) + L_k^* + \sqrt{\log(ce)/(2m)},$$

which implies the second statement of the theorem. \square

1.6.2 Selection based on a test sample

In our first application of Theorem 1.20, we assume that m independent sample pairs

$$(X'_1, Y'_1), \dots, (X'_m, Y'_m)$$

are available. This may always be achieved by simply removing m samples from the training data. Of course, this is not very attractive, but m may be small relative to n . In this case we can estimate $L(\hat{g}_k)$ by the hold-out error estimate

$$R_{n,k} = \frac{1}{m} \sum_{i=1}^m \mathbb{I}_{\hat{g}_k(X'_i) \neq Y'_i}.$$

We apply Hoeffding's inequality to show that Assumption 1 is satisfied with $c = 1$, notice that $\mathbb{E}[R_{n,k}|D_n] = L(\hat{g}_k)$, and apply Theorem 1.20 to give the following result.

COROLLARY 1.6. *Assume that the model selection algorithm is performed with the hold-out error estimate. Then*

$$\begin{aligned} & \mathbb{E}L(g_n^*) - L^* \\ & \leq \min_k \left[\mathbb{E} \left[L(\hat{g}_k) - \hat{L}_n(\hat{g}_k) \right] + \left(\inf_{g \in \mathcal{C}_k} L(g) - L^* \right) + \sqrt{\frac{\log k}{m}} \right] + \frac{1}{\sqrt{2m}}. \end{aligned}$$

In other words, the estimate achieves a nearly optimal balance between the approximation error, and the quantity

$$\mathbb{E} \left[L(\hat{g}_k) - \hat{L}_n(\hat{g}_k) \right],$$

which may be regarded as the amount of overfitting.

1.6.3 Penalization by the VC dimension

In the remaining examples we consider error estimates $R_{n,k}$ which avoid splitting the data. First recall that by the Vapnik-Chervonenkis inequality, $2\sqrt{(V_{\mathcal{C}_k} \log(n+1) + \log 2)/n}$ is an upper bound for the expected maximal deviation, within class \mathcal{C}_k , between $L(g)$ and its empirical counterpart, $\hat{L}_n(g)$. This suggests that penalizing the empirical error by this complexity term should compensate the overfitting within class \mathcal{C}_k . Thus, we introduce the error estimate

$$R_{n,k} = \hat{L}_n(\hat{g}_k) + 2\sqrt{\frac{V_{\mathcal{C}} \log(n+1) + \log 2}{n}}$$

Indeed, it is easy to show that this estimate satisfies Assumption 1. Indeed,

$$\begin{aligned}
& \mathbb{P}[L(\hat{g}_k) > R_{n,k} + \epsilon] \\
&= \mathbb{P}\left[L(\hat{g}_k) - \hat{L}_n(\hat{g}_k) > 2\sqrt{\frac{V_C \log(n+1) + \log 2}{n}} + \epsilon\right] \\
&\leq \mathbb{P}\left[\sup_{g \in \mathcal{C}_k} |L(g) - \hat{L}_n(g)| > 2\sqrt{\frac{V_C \log(n+1) + \log 2}{n}} + \epsilon\right] \\
&\leq \mathbb{P}\left[\sup_{g \in \mathcal{C}_k} |L(g) - \hat{L}_n(g)| > \mathbb{E} \sup_{g \in \mathcal{C}_k} |L(g) - \hat{L}_n(g)| + \epsilon\right] \\
&\quad (\text{by the Vapnik-Chervonenkis inequality}) \\
&\leq e^{-2n\epsilon^2} \quad (\text{by the bounded difference inequality}).
\end{aligned}$$

Therefore, satisfies Assumption 1 with $m = n$. Substituting this into Theorem 1.20 gives

$$\begin{aligned}
& \mathbb{E}L(g_n^*) - L^* \\
&\leq \min_k \left[2\sqrt{\frac{V_{C_k} \log(n+1) + \log 2}{n}} + \left(\inf_{g \in \mathcal{C}_k} L(g) - L^* \right) + \sqrt{\frac{\log k}{n}} \right] + \sqrt{\frac{1}{2n}}.
\end{aligned}$$

Thus, structural risk minimization finds the best trade-off between the approximation error and a distribution-free upper bound on the estimation error.

1.6.4 Penalization by maximum discrepancy

In this section we propose a data-dependent way of computing the penalties with improved performance guarantees. Assume, for simplicity, that n is even, divide the data into two equal halves, and define, for each predictor f , the empirical loss on the two parts by

$$\hat{L}_n^{(1)}(g) = \frac{2}{n} \sum_{i=1}^{n/2} \mathbb{I}_{g(X_i) \neq Y_i}$$

and

$$\hat{L}_n^{(2)}(g) = \frac{2}{n} \sum_{i=n/2+1}^n \mathbb{I}_{g(X_i) \neq Y_i}.$$

Define the error estimate $R_{n,k}$ by

$$R_{n,k} = \hat{L}_n(\hat{g}_k) + \max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)).$$

Observe that the maximum discrepancy $\max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g))$ may be computed using the following simple trick: first flip the labels of the first half of the data, thus obtaining the modified data set $D'_n = (X'_1, Y'_1), \dots, (X'_n, Y'_n)$ with $(X'_i, Y'_i) = (X_i, 1 - Y_i)$ for $i \leq n/2$ and $(X'_i, Y'_i) = (X_i, Y_i)$ for $i > n/2$. Next find $f_k^- \in \mathcal{C}_k$ which minimizes the empirical loss based on D'_n ,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{g(X'_i) \neq Y'_i} &= \frac{1}{2} - \frac{1}{n} \sum_{i=1}^{n/2} \mathbb{I}_{g(X_i) \neq Y_i} + \frac{1}{n} \sum_{i=n/2+1}^n \mathbb{I}_{g(X_i) \neq Y_i} \\ &= \frac{1 - \hat{L}_n^{(1)}(g) + \hat{L}_n^{(2)}(g)}{2}. \end{aligned}$$

Clearly, the function f_k^- maximizes the discrepancy. Therefore, the same algorithm that is used to compute the empirical loss minimizer \hat{g}_k may be used to find f_k^- and compute the penalty based on maximum discrepancy. This is appealing: although empirical loss minimization is often computationally difficult, the same approximate optimization algorithm can be used for both finding prediction rules and estimating appropriate penalties. In particular, if the algorithm only approximately minimizes empirical loss over the class \mathcal{C}_k because it minimizes over some proper subset of \mathcal{C}_k , the theorem is still applicable.

Theorem 1.21. *If the penalties are defined using the maximum-discrepancy error estimates, and $m = n/21$, then*

$$\begin{aligned} \mathbb{E}L(g_n^*) - L^* &\leq \min_k \left[\mathbb{E} \max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)) \right. \\ &\quad \left. + \left(\inf_{g \in \mathcal{C}_k} L(g) - L^* \right) + 4.59 \sqrt{\frac{\log k}{n}} \right] + \frac{4.70}{\sqrt{n}}. \end{aligned}$$

PROOF. Once again, we check Assumption 1 and apply Theorem 1.20. Introduce the ghost sample $(X'_1, Y'_1), \dots, (X'_n, Y'_n)$, which is independent of the data and has the same distribution. Denote the empirical loss based on this sample by $L'_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{g(X'_i) \neq Y'_i}$. The

proof is based on the simple observation that for each k ,

$$\begin{aligned}
\mathbb{E} \max_{g \in \mathcal{F}_k} (L'_n(g) - \hat{L}_n(g)) &= \frac{1}{n} \mathbb{E} \max_{g \in \mathcal{F}_k} \sum_{i=1}^n (\mathbb{I}_{g(X'_i) \neq Y'_i} - \mathbb{I}_{g(X_i) \neq Y_i}) \\
&\leq \frac{1}{n} \mathbb{E} \left(\max_{g \in \mathcal{F}_k} \sum_{i=1}^{n/2} (\mathbb{I}_{g(X'_i) \neq Y'_i} - \mathbb{I}_{g(X_i) \neq Y_i}) \right. \\
&\quad \left. + \max_{g \in \mathcal{F}_k} \sum_{i=n/2+1}^n (\mathbb{I}_{g(X'_i) \neq Y'_i} - \mathbb{I}_{g(X_i) \neq Y_i}) \right) \\
&= \frac{2}{n} \mathbb{E} \max_{g \in \mathcal{F}_k} \sum_{i=1}^{n/2} (\mathbb{I}_{g(X'_i) \neq Y'_i} - \mathbb{I}_{g(X_i) \neq Y_i}) \\
&= \mathbb{E} \max_{g \in \mathcal{F}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)). \tag{1.1}
\end{aligned}$$

The bounded difference inequality inequality (Theorem 1.8) implies

$$\mathbb{P} \left[\max_{g \in \mathcal{C}_k} (L'_n(g) - \hat{L}_n(g)) > \mathbb{E} \max_{g \in \mathcal{C}_k} (L'_n(g) - \hat{L}_n(g)) + \epsilon \right] \leq e^{-n\epsilon^2}, \tag{1.2}$$

$$\mathbb{P} \left[\max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)) < \mathbb{E} \max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)) - \epsilon \right] \leq e^{-n\epsilon^2/2} \tag{1.3}$$

and so for each k ,

$$\begin{aligned}
& \mathbb{P}[L(\hat{g}_k) > R_{n,k} + \epsilon] \\
&= \mathbb{P}\left[L(\hat{g}_k) - \hat{L}_n(\hat{g}_k) > \max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)) + \epsilon\right] \\
&\leq \mathbb{P}\left[L'_n(\hat{g}_k) - \hat{L}_n(\hat{g}_k) > \max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)) + \frac{7\epsilon}{9}\right] \\
&\quad + \mathbb{P}\left[L(\hat{g}_k) - L'_n(\hat{g}_k) > \frac{2\epsilon}{9}\right] \\
&\leq \mathbb{P}\left[L'_n(\hat{g}_k) - \hat{L}_n(\hat{g}_k) > \max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)) + \frac{7\epsilon}{9}\right] \\
&\quad + e^{-8n\epsilon^2/81} \quad (\text{by Hoeffding}) \\
&\leq \mathbb{P}\left[\max_{g \in \mathcal{C}_k} (L'_n(g) - \hat{L}_n(g)) > \max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)) + \frac{7\epsilon}{9}\right] \\
&\quad + e^{-8n\epsilon^2/81} \\
&\leq \mathbb{P}\left[\max_{g \in \mathcal{C}_k} (L'_n(g) - \hat{L}_n(g)) > \mathbb{E}\max_{g \in \mathcal{C}_k} (L'_n(g) - \hat{L}_n(g)) + \frac{\epsilon}{3}\right] \\
&\quad + \mathbb{P}\left[\max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)) < \mathbb{E}\max_{g \in \mathcal{C}_k} (\hat{L}_n^{(1)}(g) - \hat{L}_n^{(2)}(g)) - \frac{4\epsilon}{9}\right] \\
&\quad + e^{-8n\epsilon^2/81} \quad (\text{where we used (1.1)}) \\
&\leq e^{-n\epsilon^2/9} + e^{-8n\epsilon^2/81} + e^{-8n\epsilon^2/81} \quad (\text{by (1.2) and (1.3)}) \\
&< 3e^{-8n\epsilon^2/81}.
\end{aligned}$$

Thus, Assumption 1 is satisfied with $m = n/21$ and $c = 3$ and the proof is finished. \square

Bibliography

General ¹

- [1] M. Anthony and P. L. Bartlett, *Neural Network Learning: Theoretical Foundations*, Cambridge University Press, Cambridge, 1999.
- [2] N. Cristianini and J. Shawe-Taylor. *An Introduction to Support Vector Machines*. Cambridge University Press, Cambridge, UK, 2000.
- [3] L. Devroye, L. Györfi, and G. Lugosi. *A Probabilistic Theory of Pattern Recognition*. Springer-Verlag, New York, 1996.
- [4] V.N. Vapnik. *Estimation of Dependencies Based on Empirical Data*. Springer-Verlag, New York, 1982.
- [5] V.N. Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, New York, 1995.
- [6] V.N. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.
- [7] V.N. Vapnik and A.Ya. Chervonenkis. *Theory of Pattern Recognition*. Nauka, Moscow, 1974. (in Russian); German translation: *Theorie der Zeichenerkennung*, Akademie Verlag, Berlin, 1979.

Concentration for sums of independent random variables

- [8] G. Bennett. Probability inequalities for the sum of independent random variables. *Journal of the American Statistical Association*, 57:33–45, 1962.
- [9] S.N. Bernstein. *The Theory of Probabilities*. Gastehizdat Publishing House, Moscow, 1946.
- [10] H. Chernoff. A measure of asymptotic efficiency of tests of a hypothesis based on the sum of observations. *Annals of Mathematical Statistics*, 23:493–507, 1952.

¹The list of references given below contains, apart from the literature cited in the text, some of the key references in each covered topics. The list is far from being complete. Its purpose it to suggest some starting points for further reading.

- [11] T. Hagerup and C. Rüb. A guided tour of Chernoff bounds. *Information Processing Letters*, 33:305–308, 1990.
- [12] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58:13–30, 1963.
- [13] R.M. Karp. *Probabilistic Analysis of Algorithms*. Class Notes, University of California, Berkeley, 1988.
- [14] M. Okamoto. Some inequalities relating to the partial sum of binomial probabilities. *Annals of the Institute of Statistical Mathematics*, 10:29–35, 1958.

Concentration

- [15] K. Azuma. Weighted sums of certain dependent random variables. *Tohoku Mathematical Journal*, 68:357–367, 1967.
- [16] S. Boucheron, G. Lugosi, and P. Massart. A sharp concentration inequality with applications in random combinatorics and learning. *Random Structures and Algorithms*, 16:277–292, 2000.
- [17] L. Devroye. Exponential inequalities in nonparametric estimation. In G. Roussas, editor, *Nonparametric Functional Estimation and Related Topics*, pages 31–44. NATO ASI Series, Kluwer Academic Publishers, Dordrecht, 1991.
- [18] J. H. Kim. The Ramsey number $R(3, t)$ has order of magnitude $t^2 / \log t$. *Random Structures and Algorithms*, 7:173–207, 1995.
- [19] M. Ledoux. On Talagrand’s deviation inequalities for product measures. *ESAIM: Probability and Statistics*, 1, 63–87, <http://www.emath.fr/ps/>, (1996).
- [20] K. Marton. A simple proof of the blowing-up lemma. *IEEE Transactions on Information Theory*, 32:445–446, 1986.
- [21] K. Marton. Bounding \bar{d} -distance by informational divergence: a way to prove measure concentration. *Annals of Probability*, to appear:0–0, 1996.
- [22] K. Marton. A measure concentration inequality for contracting Markov chains. *Geometric and Functional Analysis*, 6:556–571, 1996. Erratum: 7:609–613, 1997.
- [23] P. Massart. About the constant in Talagrand’s concentration inequalities from empirical processes. *Annals of Probability*, 28:863–884, 2000.

- [24] C. McDiarmid. On the method of bounded differences. In *Surveys in Combinatorics 1989*, pages 148–188. Cambridge University Press, Cambridge, 1989.
- [25] W. Rhee and M. Talagrand. Martingales, inequalities, and NP-complete problems. *Mathematics of Operations Research*, 12:177–181, 1987.
- [26] J.M. Steele. An Efron-Stein inequality for nonsymmetric statistics. *Annals of Statistics*, 14:753–758, 1986.
- [27] M. Talagrand. Concentration of measure and isoperimetric inequalities in product spaces. I.H.E.S. Publications Mathématiques, 81:73–205, 1996.
- [28] M. Talagrand. New concentration inequalities in product spaces. *Invent. Math.* 126:505–563, 1996.
- [29] M. Talagrand. A new look at independence. *Annals of Probability*, 24:0–0, 1996. special invited paper.

VC theory

- [30] K. Alexander. Probability inequalities for empirical processes and a law of the iterated logarithm. *Annals of Probability*, 4:1041–1067, 1984.
- [31] M. Anthony and J. Shawe-Taylor. A result of Vapnik with applications. *Discrete Applied Mathematics*, 47:207–217, 1993.
- [32] P. Bartlett and G. Lugosi. An inequality for uniform deviations of sample averages from their means. *Statistics and Probability Letters*, 44:55–62, 1999.
- [33] L. Breiman. Bagging predictors. *Machine Learning*, 24:123–140, 1996.
- [34] Devroye, L. Bounds for the uniform deviation of empirical measures. *Journal of Multivariate Analysis*, 12:72–79, 1982.
- [35] A. Ehrenfeucht, D. Haussler, M. Kearns, and L. Valiant. A general lower bound on the number of examples needed for learning. *Information and Computation*, 82:247–261, 1989.
- [36] Y. Freund. Boosting a weak learning algorithm by majority. *Information and Computation*, 121:256–285, 1995.
- [37] E. Giné and J. Zinn. Some limit theorems for empirical processes. *Annals of Probability*, 12:929–989, 1984.

- [38] D. Haussler. Decision theoretic generalizations of the PAC model for neural net and other learning applications. *Information and Computation*, 100:78–150, 1992.
- [39] V. Koltchinskii and D. Panchenko. Empirical margin distributions and bounding the generalization error of combined classifiers, *Annals of Statistics*, 30, 2002.
- [40] M. Ledoux and M. Talagrand. *Probability in Banach Space*, Springer-Verlag, New York, 1991.
- [41] G. Lugosi. Improved upper bounds for probabilities of uniform deviations. *Statistics and Probability Letters*, 25:71–77, 1995.
- [42] D. Pollard. *Convergence of Stochastic Processes*. Springer-Verlag, New York, 1984.
- [43] R.E. Schapire, Y. Freund, P. Bartlett, and W.S. Lee. Boosting the margin: a new explanation for the effectiveness of voting methods, *Annals of Statistics*, 26:1651–1686, 1998.
- [44] R.E. Schapire. The strength of weak learnability. *Machine Learning*, 5:197–227, 1990.
- [45] M. Talagrand. Sharper bounds for Gaussian and empirical processes. *Annals of Probability*, 22:28–76, 1994.
- [46] S. Van de Geer. Estimating a regression function. *Annals of Statistics*, 18:907–924, 1990.
- [47] V.N. Vapnik. *Estimation of Dependencies Based on Empirical Data*. Springer-Verlag, New York, 1982.
- [48] V.N. Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, New York, 1995.
- [49] V.N. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.
- [50] V.N. Vapnik and A.Ya. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16:264–280, 1971.
- [51] V.N. Vapnik and A.Ya. Chervonenkis. *Theory of Pattern Recognition*. Nauka, Moscow, 1974. (in Russian); German translation: *Theorie der Zeichenerkennung*, Akademie Verlag, Berlin, 1979.
- [52] A. W. van der Vaart and J. A. Wellner. *Weak convergence and empirical processes*, Springer-Verlag, New York, 1996.

Shatter coefficients, VC dimension

- [53] P. Assouad, Sur les classes de Vapnik-Chervonenkis, *C.R. Acad. Sci. Paris*, vol. 292, Sér.I, pp. 921–924, 1981.
- [54] T. M. Cover, Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition, *IEEE Transactions on Electronic Computers*, vol. 14, pp. 326–334, 1965.
- [55] R. M. Dudley, Central limit theorems for empirical measures, *Annals of Probability*, vol. 6, pp. 899–929, 1978.
- [56] R. M. Dudley, Balls in R^k do not cut all subsets of $k+2$ points, *Advances in Mathematics*, vol. 31 (3), pp. 306–308, 1979.
- [57] P. Frankl, On the trace of finite sets, *Journal of Combinatorial Theory, Series A*, vol. 34, pp. 41–45, 1983.
- [58] D. Haussler, Sphere packing numbers for subsets of the boolean n -cube with bounded Vapnik-Chervonenkis dimension, *Journal of Combinatorial Theory, Series A*, vol. 69, pp. 217–232, 1995.
- [59] N. Sauer, On the density of families of sets, *Journal of Combinatorial Theory Series A*, vol. 13, pp. 145–147, 1972.
- [60] L. Schläffli, *Gesammelte Mathematische Abhandlungen*, Birkhäuser-Verlag, Basel, 1950.
- [61] S. Shelah, A combinatorial problem: stability and order for models and theories in infinity languages, *Pacific Journal of Mathematics*, vol. 41, pp. 247–261, 1972.
- [62] J. M. Steele, Combinatorial entropy and uniform limit laws, Ph.D. dissertation, Stanford University, Stanford, CA, 1975.
- [63] J. M. Steele, Existence of submatrices with all possible columns, *Journal of Combinatorial Theory, Series A*, vol. 28, pp. 84–88, 1978.
- [64] R. S. Wernicke and R. M. Dudley, Some special Vapnik-Chervonenkis classes, *Discrete Mathematics*, vol. 33, pp. 313–318, 1981.

Lower bounds.

- [65] A. Antos and G. Lugosi. Strong minimax lower bounds for learning. *Machine Learning*, vol.30, 31–56, 1998.

- [66] P. Assouad. Deux remarques sur l'estimation. *Comptes Rendus de l'Académie des Sciences de Paris*, 296:1021–1024, 1983.
- [67] L. Birgé. Approximation dans les espaces métriques et théorie de l'estimation. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 65:181–237, 1983.
- [68] L. Birgé. On estimating a density using Hellinger distance and some other strange facts. *Probability Theory and Related Fields*, 71:271–291, 1986.
- [69] A. Blumer, A. Ehrenfeucht, D. Haussler, and M.K. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *Journal of the ACM*, 36:929–965, 1989.
- [70] L. Devroye and G. Lugosi. Lower bounds in pattern recognition and learning. *Pattern Recognition*, 28:1011–1018, 1995.
- [71] A. Ehrenfeucht, D. Haussler, M. Kearns, and L. Valiant. A general lower bound on the number of examples needed for learning. *Information and Computation*, 82:247–261, 1989.
- [72] D. Haussler, N. Littlestone, and M. Warmuth. Predicting {0, 1}-functions on randomly drawn points. *Information and Computation*, 115:248–292, 1994.
- [73] E. Mammen, A. B. Tsybakov. Smooth discrimination analysis. *The Annals of Statistics*, 27:1808–1829, 1999.
- [74] D. Schuurmans. Characterizing rational versus exponential learning curves. In *Computational Learning Theory: Second European Conference. EuroCOLT'95*, pages 272–286. Springer Verlag, 1995.
- [75] V.N. Vapnik and A.Ya. Chervonenkis. *Theory of Pattern Recognition*. Nauka, Moscow, 1974. (in Russian); German translation: *Theorie der Zeichenerkennung*, Akademie Verlag, Berlin, 1979.

Complexity regularization

- [76] H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19:716–723, 1974.
- [77] A.R. Barron. Logically smooth density estimation. Technical Report TR 56, Department of Statistics, Stanford University, 1985.
- [78] A.R. Barron. Complexity regularization with application to artificial neural networks. In G. Roussas, editor, *Nonparametric Functional Estimation and Related Topics*, pages 561–576. NATO ASI Series, Kluwer Academic Publishers, Dordrecht, 1991.

- [79] A.R. Barron, L. Birgé, and P. Massart. Risk bounds for model selection via penalization. *Probability Theory and Related fields*, 113:301–413, 1999.
- [80] A.R. Barron and T.M. Cover. Minimum complexity density estimation. *IEEE Transactions on Information Theory*, 37:1034–1054, 1991.
- [81] P. L. Bartlett. The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network. *IEEE Transactions on Information Theory*, 44(2):525–536, March 1998.
- [82] P. Bartlett, S. Boucheron, and G. Lugosi, Model selection and error estimation. *Proceedings of the 13th Annual Conference on Computational Learning Theory*, ACM Press, pp.286–297, 2000.
- [83] L. Birgé and P. Massart. From model selection to adaptive estimation. In E. Torgersen D. Pollard and G. Yang, editors, *Festschrift for Lucien Le Cam: Research papers in Probability and Statistics*, pages 55–87. Springer, New York, 1997.
- [84] L. Birgé and P. Massart. Minimum contrast estimators on sieves: exponential bounds and rates of convergence. *Bernoulli*, 4:329–375, 1998.
- [85] Y. Freund. Self bounding learning algorithms. *Proceedings of the Eleventh Annual Conference on Computational Learning Theory*, pages 247–258, 1998.
- [86] A.R. Gallant. *Nonlinear Statistical Models*. John Wiley, New York, 1987.
- [87] S. Geman and C.R. Hwang. Nonparametric maximum likelihood estimation by the method of sieves. *Annals of Statistics*, 10:401–414, 1982.
- [88] M. Kearns, Y. Mansour, A.Y. Ng, and D. Ron. An experimental and theoretical comparison of model selection methods. In *Proceedings of the Eighth Annual ACM Workshop on Computational Learning Theory*, pages 21–30. Association for Computing Machinery, New York, 1995.
- [89] A. Krzyzak and T. Linder. Radial basis function networks and complexity regularization in function learning. *IEEE Transactions on Neural Networks*, 9:247–256, 1998.
- [90] G. Lugosi and A. Nobel. Adaptive model selection using empirical complexities. *Annals of Statistics*, vol. 27, no.6, 1999.
- [91] G. Lugosi and K. Zeger. Nonparametric estimation via empirical risk minimization. *IEEE Transactions on Information Theory*, 41:677–678, 1995.

- [92] G. Lugosi and K. Zeger. Concept learning using complexity regularization. *IEEE Transactions on Information Theory*, 42:48–54, 1996.
- [93] C.L. Mallows. Some comments on c_p . *IEEE Technometrics*, 15:661–675, 1997.
- [94] P. Massart. Some applications of concentration inequalities to statistics. *Annales de la faculté des sciences de l'université de Toulouse, Mathématiques*, série 6, **IX**:245–303, 2000.
- [95] R. Meir. Performance bounds for nonlinear time series prediction. In *Proceedings of the Tenth Annual ACM Workshop on Computational Learning Theory*, page 122–129. Association for Computing Machinery, New York, 1997.
- [96] D.S. Modha and E. Masry. Minimum complexity regression estimation with weakly dependent observations. *IEEE Transactions on Information Theory*, 42:2133–2145, 1996.
- [97] J. Rissanen. A universal prior for integers and estimation by minimum description length. *Annals of Statistics*, 11:416–431, 1983.
- [98] G. Schwarz. Estimating the dimension of a model. *Annals of Statistics*, 6:461–464, 1978.
- [99] J. Shawe-Taylor, P. L. Bartlett, R. C. Williamson, and M. Anthony. Structural risk minimization over data-dependent hierarchies. *IEEE Transactions on Information Theory*, 44(5):1926–1940, 1998.
- [100] X. Shen and W.H. Wong. Convergence rate of sieve estimates. *Annals of Statistics*, 22:580–615, 1994.
- [101] Y. Yang and A.R. Barron. Information-theoretic determination of minimax rates of convergence. *Annals of Statistics*, to appear, 1997.
- [102] Y. Yang and A.R. Barron. An asymptotic property of model selection criteria. *IEEE Transactions on Information Theory*, 44:to appear, 1998.