

0031-3203(94)00141-3

## LOWER BOUNDS IN PATTERN RECOGNITION AND LEARNING

LUC DEVROYE\*† and GÁBOR LUGOSI‡

† School of Computer Science, McGill University, 3480 University Street, Montreal, H3A 2A7 Canada

‡ Department of Mathematics, Technical University of Budapest, 1521 Stoczek u. 2, Budapest, Hungary

(Received 23 February 1994; revised 21 September 1994; received for publication 1 November 1994)

**Abstract**—Lower bounds are derived for the performance of any pattern recognition algorithm, which, using training data, selects a discrimination rule from a certain class of rules. The bounds involve the Vapnik–Chervonenkis dimension of the class, and  $L$ , the minimal error probability within the class. We provide lower bounds when  $L = 0$  (the usual assumption in Valiant’s theory of learning) and  $L > 0$ .

Learning      Nonparametric estimation      Vapnik–Chervonenkis inequality      Lower bounds  
 Pattern recognition

### 1. INTRODUCTION

In statistical pattern recognition (or classification), one is usually given a training set  $(X_1, Y_1), \dots, (X_n, Y_n)$ , which consists of  $n$  independent identically distributed  $\mathcal{X}^d \times \{0, 1\}$  valued random variables with the same distribution as  $(X, Y)$ . Denote the probability measure of  $X$  by  $\mu$ . The object is to guess  $Y$  from  $X$  and the training set. Let us formally call a given estimate (or pattern recognition rule)  $g_n(X) = g_n(X; X_1, Y_1, \dots, X_n, Y_n)$ . The best possible rule, or the Bayes rule, is the one achieving the smallest (or Bayes) probability of error,

$$L^* \stackrel{\text{def}}{=} \inf_{g: \mathcal{X}^d \rightarrow \{0, 1\}} \mathbf{P}\{g(X) \neq Y\}.$$

The object is to find rules  $g_n$  such that in a specified sense, the probability of error with  $g_n$ ,

$$L_n \stackrel{\text{def}}{=} \mathbf{P}\{g_n(X) \neq Y | X_1, Y_1, \dots, X_n, Y_n\},$$

is close to  $L^*$ .

Under the impetus of Valiant,<sup>(1)</sup> many people have recast the pattern recognition in the framework of learning. Originally this was done under two restrictions:

- $L^* = 0$ : this happens only if with probability one,  $\mathbf{P}\{Y = 1 | X\} \in \{0, 1\}$ . In pattern recognition, we speak of non-overlapping classes.

- One is interested in minimizing  $L_n$  over a given class of rules  $\mathcal{G}$ . That is, with the help of the training data, the designer picks a function from a given class of  $\{0, 1\}$ -valued functions  $\mathcal{G}$ . (In the terminology of learning theory, elements of  $\mathcal{G}$  are called concepts.) The

error with the best rule in  $\mathcal{G}$  is denoted by

$$L \stackrel{\text{def}}{=} \inf_{g \in \mathcal{G}: \mathcal{X}^d \rightarrow \{0, 1\}} \mathbf{P}\{g(X) \neq Y\}.$$

In fact, it is assumed that  $L = L^* = 0$ , that is, that the Bayes rule is in  $\mathcal{G}$ .

Later, these requirements have been relaxed. To see what the limits are that one can achieve, minimax lower bounds for the quantity

$$\sup_{(X, Y): L = 0} \mathbf{P}\{L_n - L \geq \varepsilon\}$$

are derived that are valid for all rules  $g_n$ . Needless to say, this provides us with information about the necessary sample size. An  $(\varepsilon, \delta)$  learning algorithm in the sense of Valiant<sup>(1)</sup> is one for which we may find a sample size threshold  $N(\varepsilon, \delta)$  such that for  $n \geq N(\varepsilon, \delta)$ ,

$$\sup_{(X, Y): L = 0} \mathbf{P}\{L_n - L \geq \varepsilon\} \leq \delta.$$

In this respect,  $N(\varepsilon, \delta)$  may be considered as a measure of the appropriateness of the algorithm. Blumer *et al.*<sup>(2)</sup> showed that for any algorithm,

$$N(\varepsilon, \delta) \geq C \left( \frac{1}{\varepsilon} \log \left( \frac{1}{\delta} \right) + V \right),$$

where  $C$ , is a universal constant and  $V$  is the Vapnik–Chervonenkis (or VC) dimension of  $\mathcal{G}$ , introduced by Vapnik and Chervonenkis.<sup>(3–5)</sup> We recall here that  $V$  is the largest integer  $n$  such that there exists a set  $\{x_1, \dots, x_n\} \subseteq \mathcal{X}^d$  that is shattered by  $\mathcal{G}$ . That is, for every subset  $S \subseteq \{1, \dots, n\}$ , there exists  $g \in \mathcal{G}$  such that  $g(x_i) = 1$  when  $i \in S$  and  $g(x_i) = 0$  when  $i \notin S$ . In Ehrenfeucht *et al.*<sup>(6)</sup> the lower bound was partially improved to

$$N(\varepsilon, \delta) \geq \frac{V-1}{32\varepsilon}$$

\* Author to whom correspondence should be addressed.

when  $\varepsilon \leq 1/8$  and  $\delta \leq 1/100$ . It may be combined with the previous bound.

In the first part of this note we improve this bound further in constants. More importantly, the main purpose of this note is to deal also with the case  $L > 0$ , to tie things in with the more standard pattern recognition literature. In fact, we will derive lower tail bounds as above, as well as expectation bounds for

$$\sup_{(X,Y):L\text{fixed}} E\{L_n - L\}.$$

For the  $L = 0$  case it is shown in theorem 2 that

$$\sup_{(X,Y):L=0} P\{L_n \geq \varepsilon\} \geq \frac{1}{2e\sqrt{\pi V}} \left(\frac{2e\varepsilon n}{V-1}\right)^{(V-1)/2} e^{-4n\varepsilon/(1-4\varepsilon)}.$$

Devroye and Wagner<sup>(7)</sup> showed that if  $g_n$  is a function that minimizes the empirical error

$$\sum_{i=1}^n I_{\{g(X_i) \neq Y_i\}}$$

over  $\mathcal{G}$ , and  $L = 0$ , then

$$P\{L_n \geq \varepsilon\} \leq 4 \left(\frac{2e\varepsilon n}{V}\right)^V e^{-n\varepsilon/4}.$$

( $I_{\{A\}}$  denotes the indicator of an event  $A$ .) Later this bound was improved by Blumer *et al.*<sup>(2)</sup> to

$$P\{L_n \geq \varepsilon\} \leq 2 \left(\frac{2e\varepsilon n}{V}\right)^V e^{-n\varepsilon \log 2/2}.$$

Apart from the  $\varepsilon^{(V-1)/2}$  term in the lower bound, and differences in constants, the lower bound and the upper bound have the same form.

For the case  $L > 0$ , several upper bounds for the performance of empirical error minimization were derived using Vapnik–Chervonenkis-type inequalities (see reference 8 for a survey). The best upper bounds have the form

$$c_1(n\varepsilon^2)^{c_2V} e^{-2n\varepsilon^2},$$

(see references 9–10), which are much larger than the bounds for  $L = 0$  for small  $\varepsilon$ . Among other inequalities, we show in theorem 5, that the  $\varepsilon^2$  term in the exponent is necessary. In particular, for fixed  $L \leq 1/4$

$$\sup_{(X,Y):L\text{fixed}} P\{L_n - L \geq \varepsilon\} \geq \frac{1}{4} e^{-4n\varepsilon^2/L}.$$

In general, we can conclude, that in the case  $L > 0$ , the number of samples necessary for a certain accuracy is much larger than in the usual learning theory setup, where  $L = 0$  is assumed. This phenomenon was already observed by Vapnik and Chervonenkis<sup>(11)</sup> and Simon<sup>(12)</sup> who both proved lower bounds of the type

$$\sup_{(X,Y)\text{arbitrary}} E\{L_n - L\} = \Omega\left(\sqrt{\frac{V}{n}}\right).$$

In terms of  $n$ , the order of magnitude of this lower bound is the same as those of upper bounds implied by the probability inequalities cited above. In our

theorem 3 we point out that the lower bound on the expected value of any rule depends on  $L$  as

$$\sup_{(X,Y):L\text{fixed}} E\{L_n - L\} = \Omega\left(\sqrt{\frac{LV}{n}}\right).$$

The results presented here can also be applied for classes with infinite VC dimension. For example, it is not hard to derive from theorem 1, what Blumer *et al.*<sup>(2)</sup> already pointed out, that if  $V = \infty$ , then for every  $n$  and  $g_n$ , there is a distribution with  $L = 0$  such that

$$EL_n \geq c$$

for some universal constant  $c$ . This generalizes the first theorem in reference 13, where Devroye showed a similar result if  $\mathcal{G}$  is the class of all measurable discrimination functions. Thus, when  $V = \infty$ , distribution-free nontrivial performance guarantees for  $L_n - L$  or  $L_n - L^*$  do not exist.

Other general lower bounds for  $L_n - L^*$  were also given in reference 13. For example, it is shown there that if  $L^* < 1/2$ , then for any sequence of rules  $g_n$ , and positive numbers  $a_n \rightarrow 0$ , there exists a fixed distribution such that  $EL_n \geq \min(L^* + a_n, 1/2)$  along a subsequence, that is, the rate of convergence to the Bayes-risk can be arbitrarily slow for some distributions. The difference with the minimax bounds given here is that the same distribution is used for all  $n$ , whereas the bad distributions for the bounds in this paper vary with  $n$ .

We note here that all results remain valid if we allow randomization in the algorithms  $g_n$ .

### The case $L = 0$

We begin by quoting a result by Vapnik and Chervonenkis<sup>(11)</sup> and Haussler *et al.*<sup>(14)</sup>

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

$$\sup_{(X,Y) \in \chi} EL_n \geq \frac{V-1}{2en} \left(1 - \frac{1}{n-1}\right).$$

We now turn to the probability bound. Our bound improves over the best bound that we are aware of thus far, as given in theorem 1 and corollary 5 of Ehrenfeucht *et al.*<sup>(6)</sup> In the case of  $N(\varepsilon, \delta)$ , the sample size needed for  $(\varepsilon, \delta)$  learning, the coefficient is improved by a factor of  $8/3$ .

**Theorem 2.** Let  $\mathcal{G}$  be a class of discrimination functions with VC dimension  $V \geq 2$ . Let  $\chi$  be the set of all random variables  $(X, Y)$  for which  $L = 0$ . Assume  $\varepsilon \leq 1/4$ . Define  $v = [(V - 1)/2]$ , and assume  $n \geq v$ . Then for every discrimination rule  $g_n$  based upon  $X_1, Y_1, \dots, X_n, Y_n$ ,

$$\sup_{(X,Y) \in \mathcal{X}} \mathbf{P}\{L_n \geq \varepsilon\} \geq \frac{1}{2e\sqrt{2\pi v}} \left(\frac{2ne\varepsilon}{V-1}\right)^{(V-1)/2} e^{-4ne\varepsilon/(1-4\varepsilon)}.$$

In particular, when  $\varepsilon \leq 1/8$  and

$$\log\left(\frac{1}{\delta}\right) \geq \left(\frac{4v}{e}\right) (2e\sqrt{2\pi v})^{1/v},$$

then

$$N(\varepsilon, \delta) \geq \frac{1}{8\varepsilon} \log\left(\frac{1}{\delta}\right).$$

If on the other hand  $n \geq 15$  and  $n \leq (V-1)/(12\varepsilon)$ , then

$$\sup_{(X,Y) \in \mathcal{X}} \mathbf{P}\{L_n \geq \varepsilon\} \geq \frac{1}{20}.$$

Finally, for  $\delta \leq 1/20$ , and  $\varepsilon < 1/2$ ,

$$N(\varepsilon, \delta) \geq \frac{V-1}{12\varepsilon}.$$

*Proof.* The idea is to construct a family  $\mathcal{F}$  of  $2^{V-1}$  distributions within the distributions with  $L=0$  as follows: first find points  $x_1, \dots, x_V$  that are shattered by  $\mathcal{G}$ . A member in  $\mathcal{F}$  is described by  $V-1$  bits,  $\theta_1, \dots, \theta_{V-1}$ . For convenience, this is represented as a bit vector  $\theta$ . We write  $\theta_{i-}$  and  $\theta_{i+}$  for the vector  $\theta$  in which the  $i$ th bit is set to 1 and 0, respectively. Assume  $V-1 \leq n$ . For a particular bit vector, we let  $X = x_i (i < V)$  with probability  $p$  each, while  $X = x_V$  with probability  $1-p(V-1)$ . Then set  $Y = f_\theta(X)$ , where  $f_\theta$  is defined as follows:

$$f_\theta(x) = \begin{cases} \theta_i & \text{if } x = x_i, i < V; \\ 0 & \text{if } x = x_V. \end{cases}$$

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

Observe that

$$L_n \geq p \sum_{i=1}^{V-1} I_{[g_n(x_i, X_1, Y_1, \dots, X_n, Y_n) \neq \theta_i]}.$$

Using this, for given  $\theta$ ,

$$\begin{aligned} & \mathbf{P}\{L_n \geq \varepsilon | X_1, Y_1, \dots, X_n, Y_n\} \\ & \geq \mathbf{P}\left\{p \sum_{i=1}^{V-1} I_{[g_n(x_i, X_1, Y_1, \dots, X_n, Y_n) \neq \theta_i]} \geq \varepsilon \mid X_1, Y_1, \dots, X_n, Y_n\right\}. \end{aligned}$$

This probability is either zero or one, as the event is deterministic. We now randomize and replace  $\theta$  and  $\Theta$ . For fixed  $X_1, \dots, X_n$ , we denote by  $J$  the collection  $\{j: 1 \leq j \leq V-1, \cap_{i=1}^n [X_i \neq x_j]\}$ . This is the collection of empty cells  $x_i$ . We bound our probability from below by summing over  $J$  only:

$$\begin{aligned} & \mathbf{P}\{L_n \geq \varepsilon | X_1, Y_1, \dots, X_n, Y_n\} \\ & \geq \mathbf{P}\left\{p \sum_{i \in J} I_{[g_n(x_i, \Theta, \dots)]} \geq \varepsilon \mid X_1, Y_1, \dots, X_n, Y_n\right\} \end{aligned}$$

where  $g_{ni}$  is shorthand for  $g_n(x_i, X_1, \dots, Y_n)$ . Conditionally, these are fixed members from  $\{0, 1\}$ . The  $\Theta_i$ s with  $i \in J$  constitute independent Bernoulli  $(1/2)$  random variables. Importantly, their values do not alter the  $g_{ni}$ s (this cannot be said for  $\Theta_i$  when  $i \notin J$ ). Thus, our lower bound is equal to

$$\mathbf{P}\{p \text{ Binomial}(|J|, 1/2) \geq \varepsilon | |J|\}.$$

We summarize:

$$\begin{aligned} & \sup_{(X,Y) \in \mathcal{X}} \mathbf{P}\{L_n \geq \varepsilon | X_1, Y_1, \dots, X_n, Y_n\} \\ & \geq \sup_{(X,Y) \in \mathcal{X}} \mathbf{P}\{L_n \geq \varepsilon | X_1, Y_1, \dots, X_n, Y_n\} \\ & \geq \sup_{\theta} \mathbf{P}\{L_n \geq \varepsilon | X_1, Y_1, \dots, X_n, Y_n\} \\ & \geq \mathbf{E}\{\mathbf{P}\{L_n \geq \varepsilon | X_1, Y_1, \dots, X_n, Y_n\}\} \\ & \geq \mathbf{E}\{\mathbf{P}\{p \text{ Binomial}(|J|, 1/2) \geq \varepsilon | |J|\}\}. \end{aligned}$$

As we are dealing with a symmetric binomial, it is easy to see that the last expression in the chain is at least equal to

$$\frac{1}{2} \mathbf{P}\{|J| \geq 2\varepsilon/p\}.$$

Assume that  $\varepsilon < 1/2$ . By the pigeonhole principle,  $|J| \geq 2\varepsilon/p$  if the number of points  $X_i, 1 \leq i \leq n$ , that are not equal to  $x_V$  does not exceed  $V-1-2\varepsilon/p$ . Therefore, we have a further lower bound:

$$\begin{aligned} \frac{1}{2} \mathbf{P}\{|J| \geq 2\varepsilon/p\} & \geq \frac{1}{2} \mathbf{P}\{\text{Binomial}(n, (V-1)p) \\ & \leq V-1-2\varepsilon/p\}. \end{aligned}$$

We consider two choice for  $p$ .

*Choice A.* Take  $p = 1/n$ , and assume  $12n\varepsilon \leq V-1, \varepsilon < 1/2$ . Note that for  $n \geq 15$

$$\mathbf{E}|J| = (V-1)(1-p)^n \geq \frac{V-1}{e} \left(1 - \frac{1}{n}\right) \geq \frac{V-1}{3}.$$

Also since  $0 \leq |J| \leq V-1$ , we have  $\text{Var}|J| \leq (V-1)^2/4$ . By the Chebyshev–Cantelli inequality,

$$\begin{aligned} & (1/2)\mathbf{P}\{|J| \geq 2n\varepsilon\} \\ & = (1/2)(1 - \mathbf{P}\{|J| < 2n\varepsilon\}) \\ & \geq (1/2)(1 - \mathbf{P}\{|J| < (V-1)/6\}) \\ & = (1/2)(1 - \mathbf{P}\{|J| - \mathbf{E}|J| \leq (V-1)/6 - \mathbf{E}|J|\}) \\ & \geq (1/2)(1 - \mathbf{P}\{|J| - \mathbf{E}|J| \leq -(V-1)/6\}) \\ & \geq (1/2)\left(1 - \frac{\text{Var}|J|}{\text{Var}|J| + (V-1)^2/36}\right) \\ & \geq (1/2)\left(1 - \frac{(V-1)^2/4}{(V-1)^2/4 + (V-1)^2/36}\right) \\ & = \frac{1}{20}. \end{aligned}$$

This proves the second inequality for  $\sup \mathbf{P}\{L_n \geq \varepsilon\}$ .

*Choice B.* Take  $p = 2\varepsilon/v$  and assume  $\varepsilon \leq 1/4$ . Assume  $n \geq v$ . Then the lower bound is

$$\begin{aligned} & \frac{1}{2} \mathbf{P}\{\text{Binomial}(n, 4\varepsilon) \leq v\} \\ & \geq \frac{1}{2} \binom{n}{v} (4\varepsilon)^v (1 - 4\varepsilon)^{n-v} \\ & \geq \frac{1}{2} \frac{1}{e\sqrt{2\pi v}} \left( \frac{4e\varepsilon(n-v+1)}{v(1-4\varepsilon)} \right)^v (1 - 4\varepsilon)^n \\ & \left( \text{since } \binom{n}{v} \geq \left( \frac{(n-v+1)e}{v} \right)^v \right. \\ & \times \frac{1}{e\sqrt{2\pi v}} \text{ by Stirling's formula} \left. \right) \\ & \geq \frac{1}{2} \frac{1}{e\sqrt{2\pi v}} \left( \frac{4e\varepsilon(n-v+1)}{v} \right)^v (1 - 4\varepsilon)^n \end{aligned}$$

becomes smaller as  $L$  decreases, as should be expected. The largest sample sizes are needed when  $L$  is close to  $1/2$ . (Note that for  $L = 1/2$ ,  $N(\varepsilon, \delta) = 0$ , since any random decision will give  $1/2$  error probability.) When  $L$  is very small, we provide an  $\Omega(1/n)$  lower bound, just as for the case  $L = 0$ . The constants in the bounds may be tightened at the expense of more complicated expressions.

*Theorem 3.* Let  $\mathcal{G}$  be a class of discrimination functions with VC dimension  $V \geq 2$ . Assume that  $n \geq 8(V-1)$ . Let  $\chi$  be the set of all random variables  $(X, Y)$  for which for fixed  $L \in (0, 1/2)$ ,

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

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

$$\sup_{(X, Y) \in \chi} \mathbf{E}(L_n - L) \geq \begin{cases} \sqrt{\frac{L(V-1)}{8n}} (1 - 2/n)^{2n} & \text{if } n \geq \frac{V-1}{2L} \max(4, 1/(1-2L)^2); \\ \frac{2(V-1)}{n} (1 - 8/n)^{2n} & \text{if } n \leq 2(V-1)/L \text{ (this implies } L \leq 1/4). \end{cases}$$

$$\begin{aligned} & \geq \frac{1}{2} \frac{1}{e\sqrt{2\pi v}} \left( \frac{4en\varepsilon}{v} \right)^v (1 - 4\varepsilon)^n \left( 1 - \frac{v-1}{n} \right)^v \\ & \geq \frac{1}{2} \frac{1}{e\sqrt{2\pi v}} \left( \frac{2en\varepsilon}{v} \right)^v e^{-4n\varepsilon/(1-4\varepsilon)} \quad (\text{since } n \geq 2(v-1)) \\ & (\text{use } 1-x \geq \exp(-x/(1-x))) \\ & \geq \frac{1}{2} \frac{1}{e\sqrt{2\pi v}} \left( \frac{2en\varepsilon}{v} \right)^v e^{-8n\varepsilon} \quad (\text{since } \varepsilon \leq 1/8) \\ & \geq \frac{(8\varepsilon)^v}{\log^v(1/\delta)} n^v e^{-8n\varepsilon} \\ & \left( \text{since we assume } \log \left( \frac{1}{\delta} \right) \geq \left( \frac{4v}{e} \right) (2e\sqrt{2\pi v})^{1/v} \right). \end{aligned}$$

The function  $n^v e^{-8n\varepsilon}$  varies unimodally in  $n$ , and achieves a peak at  $n = v/(8\varepsilon)$ . For  $n$  below this threshold, by monotonicity, we apply the bound at  $n = v/(8\varepsilon)$ . It is easy to verify that the value of the bound at  $v/(8\varepsilon)$  is always at least  $\delta$ . If on the other hand,  $(1/8\varepsilon) \log(1/\delta) \geq n \geq v/(8\varepsilon)$ , the lower bound achieves its minimal value at  $(1/8\varepsilon) \log(1/\delta)$ , and the value there is  $\delta$ . This concludes the proof.

*The case  $L > 0$*

In this section, we consider both expectation and probability bounds when  $L > 0$ . The bounds involve  $n$ ,  $V$  and  $L$  jointly. The minimax lower bound below is valid for any discrimination rule, and depends upon  $n$  as  $\sqrt{L(V-1)/n}$ . As a function of  $n$ , this decrease as in the central limit theorem. Interestingly, the lower bound

*Proof.* Again we consider the finite family  $\mathcal{F}$  from the previous section. The notation  $\theta$  and  $\Theta$  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

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

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  $\theta$  is the one which sets

$$f_\theta(x) = \begin{cases} 1 & \text{if } x = x_i, i < V, \theta_i = 1; \\ 0 & \text{otherwise.} \end{cases}$$

Also, observe that

$$L = (V-1)p(1/2 - c). \tag{1}$$

We may then write, for fixed  $\theta$ ,

$$L_n - L \geq \sum_{i=1}^{V-1} 2pc I_{\{g_n(x_i, X_1, Y_1, \dots, X_n, Y_n) = 1 - f_\theta(x_i)\}}$$

It is sometimes convenient to make the dependence of  $g_n$  upon  $\theta$  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  $\theta_i$ . The proof below is based upon Hellinger distances, and its methodology is essentially due to Assouad.<sup>(15)</sup> We replace  $\theta$  by a

uniformly distributed random  $\Theta$  over  $\{0, 1\}^{V-1}$ . Thus,

$$\begin{aligned} \sup_{(X, Y) \in \mathcal{F}} \mathbf{E}\{L_n - L\} &= \sup_{\theta} \mathbf{E}\{L_n - L\} \\ &\geq \mathbf{E}\{L_n - L\} \quad (\text{with random } \Theta) \\ &\geq \sum_{i=1}^{V-1} 2pc \mathbf{E} I_{\{g_n(x_i, x_1, \dots, x_{n-1}) = f_n(x_i)\}} \end{aligned}$$

Fix  $i < V$  and call the  $i$ th summand in the last expression  $E_i$ . Introduce the notation

$$\begin{aligned} p_{\theta}(x'_1, \dots, x'_n, y_1, \dots, y_n) \\ = \mathbf{P}\{\cap_{j=1}^n [X_j = x'_j, Y_j = y_j] | \Theta = \theta\}. \end{aligned}$$

Clearly, this may be written as a Cartesian product:

$$p_{\theta}(x'_1, \dots, x'_n, y_1, \dots, y_n) = \prod_{j=1}^n p_{\theta}(x'_j, y_j),$$

where  $p_{\theta}(x'_j, y_j) = \mathbf{P}\{X = x'_j, Y = y_j | \Theta = \theta\}$ . Thus,

yields the following:

$$\begin{aligned} E_i &\geq \frac{c^p}{2} 2^{-(V-1)} \\ &\quad \times \sum_{\theta} \left( \sum_{(x, y)} \sqrt{p_{\theta_i}(x, y) p_{\theta_i}(x, y)} \right)^{2n} \\ &= \frac{c^p}{2} 2^{-(V-1)} \sum_{\theta} \left( \sum_{\substack{(x, y) \\ x \neq x_i}} p_{\theta}(x, y) + 2p\sqrt{1/4 - c^2} \right)^{2n} \\ &= \frac{c^p}{2} 2^{-(V-1)} \sum_{\theta} \left( \sum_{(x, y)} p_{\theta}(x, y) + 2p\sqrt{1/4 - c^2} \right. \\ &\quad \left. - p_{\theta}(x_i, 1) - p_{\theta}(x_i, 0) \right)^{2n} \\ &= \frac{c^p}{2} 2^{-(V-1)} \sum_{\theta} (1 + p\sqrt{1 - 4c^2} - p)^{2n} \\ &= \frac{c^p}{2} (1 + p\sqrt{1 - 4c^2} - p)^{2n}. \end{aligned}$$

$$\begin{aligned} E_i &= 2pc 2^{-(V-1)} \sum_{\substack{(x'_1, \dots, x'_n, y_1, \dots, y_n) \\ \in \{0, 1\}^n \times \{0, 1\}^n}} \sum_{\theta} I_{\{g_n(x_i, x'_1, \dots, x'_{n-1}) = f_n(x_i)\}} \prod_{j=1}^n p_{\theta}(x'_j, y_j) \\ &= 2pc 2^{-(V-1)} \sum_{(x'_1, \dots, x'_n, y_1, \dots, y_n)} \sum_{\theta} \frac{1}{2} \left\{ I_{\{g_n(x_i, x'_1, \dots, x'_{n-1}) = 1\}} \prod_{j=1}^n p_{\theta_i}(x'_j, y_j) + I_{\{g_n(x_i, x'_1, \dots, x'_{n-1}) = 0\}} \prod_{j=1}^n p_{\theta_i}(x'_j, y_j) \right\} \\ &\geq 2pc 2^{-(V-1)} \sum_{(x'_1, \dots, x'_n, y_1, \dots, y_n)} \sum_{\theta} \min \left( \prod_{j=1}^n p_{\theta_i}(x'_j, y_j), \prod_{j=1}^n p_{\theta_i}(x'_j, y_j) \right) \\ &\quad \times \frac{1}{2} \{ I_{\{g_n(x_i, x'_1, \dots, x'_{n-1}) = 1\}} + I_{\{g_n(x_i, x'_1, \dots, x'_{n-1}) = 0\}} \} \\ &= cp 2^{-(V-1)} \sum_{(x'_1, \dots, x'_n, y_1, \dots, y_n)} \sum_{\theta} \min \left( \prod_{j=1}^n p_{\theta_i}(x'_j, y_j), \prod_{j=1}^n p_{\theta_i}(x'_j, y_j) \right) \\ &\geq \frac{c^p}{2} 2^{-(V-1)} \sum_{\theta} \left( \sum_{(x'_1, \dots, x'_n, y_1, \dots, y_n)} \sqrt{\prod_{j=1}^n p_{\theta_i}(x'_j, y_j) \times \prod_{j=1}^n p_{\theta_i}(x'_j, y_j)} \right)^2 \\ &= \frac{c^p}{2} 2^{-(V-1)} \sum_{\theta} \left( \sum_{(x, y)} \sqrt{p_{\theta_i}(x, y) p_{\theta_i}(x, y)} \right)^{2n}, \end{aligned}$$

(by the identity  $(\sum_i a_i)^n = \sum_{i_1, \dots, i_n} a_{i_1} \dots a_{i_n}$ ) where we used a discrete version of LeCam's inequality (reference 19; for example see page 7 of reference 18), which states that for positive sequences  $a_i$  and  $b_i$ , both summing to one,

$$\sum_i \min(a_i, b_i) \geq \frac{1}{2} \left( \sum_n \sqrt{a_i b_i} \right)^2.$$

We note next that for  $x = x_j, 1 \leq j \leq V, j \neq i$ ,

$$p_{\theta_i}(x, y) = p_{\theta_i}(x, y) = p_{\theta}(x, y).$$

For  $x = x_i, i < V$ , we have

$$p_{\theta_i}(x, y) p_{\theta_i}(x, y) = p^2 \left( \frac{1}{4} - c^2 \right).$$

Resubstitution in the previous chain of inequalities

As the right-hand-side does not depend upon  $i$ , the overall bound becomes

$$\begin{aligned} \sup_{(X, Y) \in \mathcal{F}} \mathbf{E}(L_n - L) &\geq \sum_{i=1}^{V-1} E_i \\ &\geq \frac{(V-1)cp}{2} (1 + p\sqrt{1 - 4c^2} - p)^{2n}. \end{aligned}$$

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

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

This 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  $2nL/(V-1) \geq 4$ , then  $c \leq 1/4$ .

With this choice for  $c$ , the lower bound is

$$\begin{aligned} \sup_{(X,Y) \in \mathcal{F}} \mathbf{E}(L_n - L) &\geq \frac{Lc}{1-2c} (1-4pc^2)^{2n} \\ &\quad (\text{since } L = (V-1)p(1/2-c) \text{ and} \\ &\quad \sqrt{1-x} - 1 \geq -x \text{ for } 0 \leq x \leq 1) \\ &\geq Lc(1-1/(n(1-2c)))^{2n} \\ &\quad (\text{by our choice of } c \text{ and the} \\ &\quad \text{expression for } L \text{ above}) \\ &\geq \sqrt{\frac{(V-1)L}{8n}} (1-2/n)^{2n} \\ &\quad (\text{since } c \leq 1/4). \end{aligned}$$

The condition  $p(V-1) \leq 1$  implies that we need to ask that  $n \geq (V-1)/(2L(1-2L)^2)$ .

Assume next that  $2nL/(V-1) \leq 4$ . Then we may put  $p = 8/n$ . Assume that  $n \geq 8(V-1)$ . This leads to a value of  $c$  determined by  $1-2c = nL/4(V-1)$ . In that case, as  $c \geq 1/4$ , the overall lower bound may be written as

$$\frac{(V-1)cp(1-p)^{2n}}{2} \geq (1-8/n)^{2n} \frac{2(V-1)}{n}.$$

This concludes the proof of theorem 3.

From the expectation bound in theorem 3, we may derive a probabilistic bound by a rather trivial argument. Unfortunately, the bound thus obtained only yields a suboptimal estimate for  $N(\varepsilon, \delta)$ .

*Theorem 4.* Let  $\mathcal{G}$  be a class of discrimination functions with VC dimension  $V \geq 2$ . Assume that  $n \geq 8(V-1)$ . Let  $\chi$  be the set of all random variables  $(X, Y)$  for which for fixed  $L \in (0, 1/2)$ ,

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

Then, for every discrimination rule  $g_n$  based upon  $X_1, Y_1, \dots, X_n, Y_n$  and any  $\varepsilon \leq A/2$ ,

$$\sup_{(X,Y) \in \chi} \mathbf{P}\{L_n - L \geq \varepsilon\} \geq \frac{A}{2-A},$$

where

$$A = \begin{cases} \sqrt{\frac{L(V-1)}{8n}} (1-2/n)^{2n} & \text{if } n \geq \frac{V-1}{2L} \max(4, 1/(1-2L)^2); \\ \frac{2(V-1)}{n} (1-8/n)^{2n} & \text{if } n \leq 2(V-1)/L. \end{cases}$$

Also,

$$N(\varepsilon, \delta) \geq \frac{L(V-1)e^{-10}}{32} \times \min\left(\frac{1}{\delta^2}, \frac{1}{\varepsilon^2}\right).$$

*Proof.* Assume that we have  $\mathbf{E}(L_n - L) \geq A$  (as in theorem 3). Then a simple bounding argument yields, for  $\varepsilon \leq A$ ,

$$\mathbf{P}\{L_n - L \geq \varepsilon\} \geq \frac{A - \varepsilon}{1 - \varepsilon}.$$

For  $\varepsilon \leq A/2$ , the lower bound is at least  $A/(2-A)$ .

For the bound on  $N(\varepsilon, \delta)$  assume that  $\mathbf{P}\{L_n - L > \varepsilon\} < \delta$ . Then clearly,  $\mathbf{E}\{L_n - L\} \leq \varepsilon + \delta$ . Thus, when  $n$  is large enough to satisfy the assumptions of theorem 3, we have

$$\sqrt{\frac{L(V-1)}{8n}} \left(1 - \frac{2}{n}\right)^{2n} \leq \varepsilon + \delta.$$

Note that

$$(1-2/n)^{2n} \geq \exp\left(-\frac{4}{1-2/n}\right) \geq e^{-5}$$

when  $n \geq 10$ . We have

$$n \geq \frac{L(V-1)}{8e^{10}(\varepsilon + \delta)^2} \geq \frac{L(V-1)}{32e^{10}} \times \min\left(\frac{1}{\delta^2}, \frac{1}{\varepsilon^2}\right).$$

It is easy to see that  $\sup_{(X,Y)} \inf_{g_n} \mathbf{P}\{L_n - L > \varepsilon\}$  is monotone decreasing in  $n$ , therefore, small values of  $n$  cannot lead to better bounds for  $N(\varepsilon, \delta)$ .

*Theorem 5.* Let  $\mathcal{G}$  be a class of discrimination functions with VC dimension  $V \geq 2$ . Let  $\chi$  be the set of all random variables  $(X, Y)$  for which for fixed  $L \in (0, 1/4)$ ,

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

Then, for every discrimination rule  $g_n$  based upon  $X_1, Y_1, \dots, X_n, Y_n$  and any  $\varepsilon \leq L$ ,

$$\sup_{(X,Y) \in \chi} \mathbf{P}\{L_n - L \geq \varepsilon\} \geq \frac{1}{4} e^{-4n\varepsilon^2/L},$$

and in particular, for  $\varepsilon \leq L \leq 1/4$ ,

$$N(\varepsilon, \delta) \geq \frac{L}{4\varepsilon^2} \log \frac{1}{4\delta}.$$

*Proof.* The argument here is similar to that in the proof of theorem 3. Using the same notation as there, it is clear that

$$\begin{aligned} \sup_{(X,Y) \in \chi} \mathbf{P}\{L_n - L \geq \varepsilon\} &\geq \mathbf{E} I_{\left\{ \sum_{i=1}^{V-1} 2^{pi} I_{\{g_n(x_i, X_1, \dots, Y_n) = 1 - f_{\theta}(x_i)\} \geq \varepsilon} \right\}} \\ &= 2^{-(V-1)} \sum_{\substack{(x'_1, \dots, x'_n, y_1, \dots, y_n) \\ \in (\{0,1\}^n \times \{0,1\})^n}} \sum_{\theta} I_{\left\{ \sum_{i=1}^{V-1} 2^{pi} I_{\{g_n(x_i, x'_1, y_1, \dots, x'_n, y_n) = 1 - f_{\theta}(x_i)\} \geq \varepsilon} \right\}} \prod_{j=1}^n p_{\theta}(x'_j, y_j). \end{aligned}$$

Now, observe, that if  $\epsilon/(2pc) \leq (V-1)/2$  (which will be called condition (\*) below), then

$$I_{[\sum_{i=1}^V 2^{pc} I_{\{\theta_{i^c}(x_1, y_1, \dots, x_n, y_n) \geq \epsilon\}}]} + I_{[\sum_{i=1}^V 2^{pc} I_{\{\theta_i(x_1, y_1, \dots, x_n, y_n) \geq \epsilon\}}]} \geq 1,$$

where  $\theta^c$  denotes the binary vector  $(1 - \theta_1, \dots, 1 - \theta_{V-1})$ , that is, the complement of  $\theta$ . Therefore, for  $\epsilon \leq pc(V-1)$ , the last expression in the lower bound above is bounded from below by

$$\frac{1}{2^{V-1}} \sum_{\substack{(x_1, \dots, x_n, y_1, \dots, y_n) \\ \in (\{0,1\}^n \times \{0,1\}^n)}} \sum_{\theta} \frac{1}{2} \min \left( \prod_{j=1}^n p_{\theta}(x'_j, y_j), \prod_{j=1}^n p_{\theta^c}(x'_j, y_j) \right) \\ \geq \frac{1}{2^{V-1}} \sum_{\theta} \left( \sum_{\substack{(x_1, \dots, x_n, y_1, \dots, y_n) \\ \in (\{0,1\}^n \times \{0,1\}^n)}} \sqrt{\prod_{j=1}^n p_{\theta}(x'_j, y_j) \times \prod_{j=1}^n p_{\theta^c}(x'_j, y_j)} \right)^2 = \frac{1}{2^{V-1}} \sum_{\theta} \left( \sum_{(x,y)} \sqrt{p_{\theta}(x, y) p_{\theta^c}(x, y)} \right)^{2n}$$

It is easy to see that for  $x = x_v$

$$p_{\theta}(x, y) = p_{\theta^c}(x, y) = \frac{1 - (V-1)p}{2}$$

and for  $x = x_i, i < V$ ,

$$p_{\theta}(x, y) p_{\theta^c}(x, y) = p^2 \left( \frac{1 - c^2}{4} \right).$$

Thus, we have the equality

$$\sum_{(x,y)} \sqrt{p_{\theta}(x, y) p_{\theta^c}(x, y)} = 1 - (V-1)p \\ + 2(V-1)p \sqrt{\frac{1 - c^2}{4}}$$

Summarizing, since  $L = p(V-1)(1/2 - c)$ , we have

$$\sup_{(X,Y) \in \mathcal{Z}} \mathbf{P}\{L_n - L \geq \epsilon\} \geq \frac{1}{4} \left( 1 - \frac{L}{1 - c} \left( 1 - \sqrt{1 - 4c^2} \right) \right)^{2n} \\ \geq \frac{1}{4} \left( 1 - \frac{L}{1 - c} 4c^2 \right)^{2n} \\ \geq \frac{1}{4} \exp \left( - \frac{16nLc^2}{1 - 2c} \left/ \left( 1 - \frac{8Lc^2}{1 - 2c} \right) \right. \right),$$

where again, we used the inequality  $1 - x \geq e^{-x/(1-x)}$ . We may choose  $c$  as  $\epsilon/(2L + 2\epsilon)$ . It is easy to verify that condition (\*) holds. Also,  $p(V-1) \leq 1$ . From the condition  $L \geq \epsilon$  we deduce that  $c \leq 1/4$ . The exponent in the expression above may be bounded as

$$\frac{16nLc^2}{1 - 2c} = \frac{16nLc^2}{1 - 2c - 8Lc^2} \\ = \frac{4n\epsilon^2}{L + \epsilon} \\ = \frac{2\epsilon^2}{L + \epsilon}$$

(by substituting  $c = \epsilon/(2L + 2\epsilon)$ )

$$\leq \frac{4n\epsilon^2}{L}$$

Thus,

$$\sup_{(X,Y) \in \mathcal{Z}} \mathbf{P}\{L_n - L \geq \epsilon\} \geq \frac{1}{4} \exp(-4n\epsilon^2/L),$$

as desired. Setting this bound equal to  $\delta$  provides the bound on  $N(\epsilon, \delta)$ .

Theorems 4 and 5 may of course be combined. They show that  $N(\epsilon, \delta)$  is bounded from below by terms like  $(1/\epsilon^2) \log(1/\delta)$  (independent of  $V$ ) and  $(V-1)/\epsilon^2$ . As  $\delta$  is typically small, the main term is thus not influenced by the VC dimension. This is the same phenomenon as in the case  $L = 0$ .

*Acknowledgements*—We thank Dr Simon for sending us his 1993 paper on lower bounds. The first author's research was sponsored by NSERC grant A3456 and FCAR grant 90-ER-0291.

REFERENCES

1. L. G. Valiant, A theory of the learnable, *Commun. ACM* **27**, 1134–1142 (1984).
2. A. Blumer, A. Ehrenfeucht, D. Haussler and M. Warmuth, Learnability and the Vapnik–Chervonenkis dimension, *J. ACM* **36**, 0–0 (1989).
3. V. N. Vapnik and A. Ya. Chervonenkis, On the uniform convergence of relative frequencies of events to their probabilities, *Theory Prob. Appl.* **16**, 264–280 (1971).
4. V. N. Vapnik and A. Ya. Chervonenkis, Theory of uniform convergence of frequencies of events to their probabilities and problems of search for an optimal solution from empirical data, *Automat. Remote Cont.* **32**, 207–217 (1971).
5. V. N. Vapnik and A. Ya. Chervonenkis, Necessary and sufficient conditions for the uniform convergence of means to their expectations, *Theory Probab. Appl.* **26**, 532–553 (1981).
6. A. Ehrenfeucht, D. Haussler, M. Kearns and L. Valiant, A general lower bound on the number of examples needed for learning, *Inform. Comput.* **82**, 247–261 (1989).
7. L. Devroye and T. J. Wagner, Nonparametric discrimination and density estimation, Technical report 183, Electronics Research Center, University of Texas (1976).
8. L. Devroye, Automatic pattern recognition: a study of the probability of error, *IEEE Trans. Pattern Anal. Mach. Intell.* **10**, 530–543 (1988).
9. K. Alexander, Probability inequalities for empirical processes and a law of the iterated logarithm, *Ann. Probab.* **12**, 1041–1067 (1984).
10. M. Talagrand, Sharper bounds for Gaussian and empirical processes, *Ann. Probab.* **22**, 28–76 (1994).
11. V. N. Vapnik and A. Ya. Chervonenkis, *Theory of Pattern Recognition*, Nauka, Moscow (1974).
12. H. U. Simon, General lower bounds on the number of examples needed for learning probabilistic concepts, *Proc.*

- 6th Annual Workshop on Computational Learning Theory, pp. 402–412 (1993).
13. L. Devroye, Any discrimination rule can have an arbitrarily bad probability of error for finite sample size, *IEEE Trans. Pattern Anal. Mach. Intell.* **4**, 154–157 (1982).
  14. D. Haussler, N. Littlestone and M. Warmuth, Predicting  $\{0, 1\}$  functions from randomly drawn points *Proc. 29th IEEE Symp. on the Foundations of Computer Science*, 100–109 (1988).
  15. P. Assouad, Deux remarques sur l'estimation, *C.R.Acad. Sci. Paris*, **296**, 1021–1024 (1983).
  16. L. Birgé, Approximation dans les espaces métriques et théorie de l'estimation, *Z. Wahrscheinlich. Ver. Gebiete* **65**, 181–237 (1983).
  17. L. Birgé, On estimating a density using Hellinger distance and some other strange facts, *Probab. Theory Related Fields*, **71**, 271–291 (1986).
  18. L. Devroye, *A Course in Density Estimation*. Birkhäuser, Boston (1987).
  19. L. LeCam, Convergence of estimates under dimensionality restrictions, *Ann. Stat.* **1**, 38–53 (1973).

**About the Author**— LUC DEVROYE was born in Belgium on 6 August, 1948. He obtained his doctoral degree from the University of Texas at Austin in 1976 and joined McGill University in 1977, where he is Professor of Computer Science. He is interested in probability theory, probabilistic analysis of data structures and algorithms, random search, pattern recognition, nonparametric estimation, random variate generation, and font design.

**About the Author**— GÁBOR LUGOSI was born on 13 July, 1964 in Budapest, Hungary. He is an associate professor at the Department of Mathematics and Computer Science, Technical University of Budapest. He received his degree in electrical engineering from the Technical University of Budapest in 1987, and his Ph.D. from the Hungarian Academy of Sciences in 1991. His main research interests include pattern recognition, information theory, and nonparametric statistics.