

## Lower Bounds for Bayes Error Estimation

András Antos, Luc Devroye, and  
László Györfi, *Fellow, IEEE*

**Abstract**—We give a short proof of the following result. Let  $(X, Y)$  be any distribution on  $\mathcal{N} \times \{0, 1\}$ , and let  $(X_1, Y_1), \dots, (X_n, Y_n)$  be an i.i.d. sample drawn from this distribution. In discrimination, the Bayes error  $L^* = \inf_g \mathbf{P}\{g(X) \neq Y\}$  is of crucial importance. Here we show that without further conditions on the distribution of  $(X, Y)$ , no rate-of-convergence results can be obtained. Let  $\phi_n(X_1, Y_1, \dots, X_n, Y_n)$  be an estimate of the Bayes error, and let  $\{\phi_n(\cdot)\}$  be a sequence of such estimates. For any sequence  $\{a_n\}$  of positive numbers converging to zero, a distribution of  $(X, Y)$  may be found such that  $\mathbf{E}\{|L^* - \phi_n(X_1, Y_1, \dots, X_n, Y_n)|\} \geq a_n$  infinitely often.

**Index Terms**—Discrimination, statistical pattern recognition, nonparametric estimation, Bayes error, lower bounds, rates of convergence.

### 1 INTRODUCTION

The pattern recognition problem may be formulated as follows: we are given  $n$  i.i.d. observations  $D_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ , drawn from the common unknown distribution of  $(X, Y)$  on  $\mathbb{R}^d \times \{0, 1\}$ . Given  $X$ , one must estimate  $Y$  as best as possible by a function  $g_n(X)$  of  $X$  and the observations. The best one can hope for is to make an error equal to the Bayes error,  $L^*$ :

$$L_n \stackrel{\text{def}}{=} \mathbf{P}\{g_n(X) \neq Y | D_n\} \geq L^* \stackrel{\text{def}}{=} \inf_{g: \mathbb{R}^d \rightarrow \{0, 1\}} \mathbf{P}\{g(X) \neq Y\}.$$

It is thus of great importance to be able to estimate  $L^*$  accurately, even before pattern recognition is attempted. Also, a comparison of estimates of  $L_n$  and  $L^*$  gives us an idea how much room is left for improvement.

In a first group of methods,  $L^*$  is estimated by an estimate  $\hat{L}_n$  of the error probability  $L_n$  of some consistent classification rule  $g_n$ . As such, this problem has been attempted by Fukunaga and Kessel [9], Chen and Fu [2], Fukunaga and Hummels [8], and Garnett and Yau [10], to cite just the early contributions. Concerning the error estimation of specific classification rules see Chapter 10 in McLachlan [11]. Clearly, if the estimate  $\hat{L}_n$  we use is consistent in the sense that  $\hat{L}_n - L_n \rightarrow 0$  with probability one as  $n \rightarrow \infty$ , and the rule is strongly consistent, then  $\hat{L}_n \rightarrow L^*$  with probability one. In other words, we have a consistent estimate of the Bayes error probability. The problem is that even though for many classifiers,  $\hat{L}_n - L_n$  can be guaranteed to converge to zero rapidly, regardless what the distribution of  $(X, Y)$  is (see Chapters 8, 23, 24, and 31 of Devroye et al. [7]), in view of Cover [3] and Devroye [4], the rate of convergence of  $L_n$  to  $L^*$  using such a method may be arbitrarily slow. Thus, we cannot expect a good performance for all distributions from such a method. The question thus is whether it is possible to come up with another method of estimating  $L^*$  (by  $\phi_n(X_1, Y_1, \dots, X_n, Y_n)$ ) such that the difference

$\phi_n(X_1, Y_1, \dots, X_n, Y_n) - L^*$  converges to zero rapidly for all distributions. Unfortunately, there is no method that guarantees a certain finite sample performance for all distributions. This disappointing fact is reflected in the following negative result (Theorem 8.5 of Devroye et al. [7]).

**Theorem 1.** *For every  $n$ , for any estimate  $\phi_n(X_1, Y_1, \dots, X_n, Y_n)$  of the Bayes error probability  $L^*$ , and for every  $\epsilon > 0$ , there exists a distribution of  $(X, Y)$ , such that*

$$\mathbf{E}\{|\phi_n(X_1, Y_1, \dots, X_n, Y_n) - L^*|\} \geq 1/4 - \epsilon.$$

The counterexamples in Theorem 1 vary with  $n$ , so it may still be possible that for every fixed distribution for  $(X, Y)$ , there exists a universal rate of convergence to zero for

$$\mathbf{E}\{|\phi_n(X_1, Y_1, \dots, X_n, Y_n) - L^*|\}.$$

The purpose of this note is to show that this too is impossible. We show the following:

**Theorem 2.** *For any sequence  $\{a_n\}$  of positive numbers converging to zero, a distribution of  $(X, Y)$  on  $\{1, 2, 3, \dots\} \times \{0, 1\}$  may be found such that*

$$\mathbf{E}\{|\phi_n(X_1, Y_1, \dots, X_n, Y_n) - L^*|\} \geq a_n$$

*infinitely often.*

We note that for the  $L_1$  error in density estimation, similar global lower bounds were obtained by Devroye [5], [6] and Birgé [1]. We also note that the phrase “infinitely often” cannot be dropped from Theorem 2. Indeed, there exist deterministic sequences  $b_n$  with  $|b_n - L^*| \leq c/n$  infinitely often for some constant  $c$ : just consider the dyadic sequence

$$0/2^0, 1/2^0, 0/2^1, 1/2^1, 2/2^1, 0/2^2, 1/2^2, 2/2^2, 3/2^2, 4/2^2, \dots$$

With  $\phi_n \equiv b_n$ , we thus obtain a very good estimate along an (unknown) subsequence.

### 2 PROOF OF THEOREM 2

Given  $\{a_n\}$ , we find a sequence of positive integers  $\ell_n$  with a given property to be specified later. Then, we partition the positive integers into consecutive blocks of cardinality  $\ell_1, \ell_2, \ell_3, \dots$ . Let  $z = (z_1, z_2, \dots)$  be a vector assigning a bit to each integer, and let  $u = (u_1, u_2, \dots)$  be a vector assigning a bit to each block. Then the distribution of  $(X, Y)$  is described constructively as follows: first a block  $B$  is drawn from the geometric distribution:

$$\mathbf{P}\{B = i\} = \frac{1}{2^i}, \quad i \geq 1.$$

Then  $X$  is drawn uniformly over the  $\ell_B$  integers in that block. If  $u_B = 0$ , then  $Y = z_X$ , while if  $u_B = 1$ ,  $Y$  is Bernoulli (1/2). For this distribution, it is easy to verify that

$$L^* = L^*(u) = \sum_{i=1}^{\infty} \frac{u_i}{2^{i+1}}$$

as the only problem blocks are those with  $u_i = 1$ , where locally, the Bayes error conditioned on  $X \in \text{block } i$  is 1/2. Note in particular that  $L^*$  depends upon  $u$  only.

Assume that all samples have the same common  $X_1, X_2, \dots, X_n$  components, consisting of i.i.d. observations drawn from the distribution of  $X$  (which is the same for all  $(u, z)$ ). Then let  $W = (W_1, \dots, W_n)$  be a bit vector consisting of i.i.d. Bernoulli (1/2) random bits. Then, define

- A. Antos is with the Computer and Automation Research Institute of the Hungarian Academy of Sciences, Lágymányosi u. 11, Budapest, Hungary H-1518.
- L. Devroye is with the School of Computer Science, McGill University, Montreal, Canada H3A 2K6.
- L. Györfi is with the Department of Computer Science and Information Theory, Technical University of Budapest, Stoczek u. 2, Budapest, Hungary H-1521. E-mail: györfi@inf.bme.hu.

Manuscript received 28 Apr. 1998; revised 23 Mar. 1999.

Recommended for acceptance by A. Webb.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number 107667.

$$Y_i = \begin{cases} z_{X_i} & \text{if } u_b = 0 \text{ where } b \text{ is the block number of } X_i \\ W_i & \text{otherwise.} \end{cases}$$

Put

$$D_n = D_n(u, z) = \{(X_1, Y_1), \dots, (X_n, Y_n)\}.$$

In the proof we use randomization such that  $(u, z)$  is replaced by the independent random sequences  $(U, Z)$ , where  $U = (U_1, U_2, \dots)$  and  $Z = (Z_1, Z_2, \dots)$  are i.i.d. Bernoulli  $(1/2)$  sequences.

Let  $U_+^k$  and  $U_-^k$  denote the vector  $U$ , with the difference that  $U_+^k$  forces the  $k$ th bit to be 1 and  $U_-^k$  forces the  $k$ th bit to be 0. The event that all  $X_1, \dots, X_n$  belong to block  $k$  are disjoint is denoted by  $T_k$ .

Introduce the notation

$$R_n(u, z) = \mathbf{E} \{ |\phi_n(D_n(u, z)) - L^*(u)| \},$$

$$\bar{R}_n(u, z) = \sup_{m > n} \frac{R_m(u, z)}{a_m}$$

and

$$A_n = \{(u, z) : \bar{R}_n(u, z) \leq 1\}.$$

Observe that given  $X_1, \dots, X_n$  and  $U$  and the event  $T_k$ , the distributions of  $(D_n(U_-^k, Z), I_{\{(U_+^k, Z) \in A_n\}})$  and  $(D_n(U_+^k, Z), I_{\{(U_+^k, Z) \in A_n\}})$  are equal. In order to see this decompose  $D_n(U_+^k, Z)$  and  $D_n(U_-^k, Z)$  as follows: Let  $X_{i_1}, \dots, X_{i_M}$  be the subsample of  $X_1, \dots, X_n$  such that  $X_{i_j}$ s fall into block  $k$ .

Put

$$E_n(U, Z) = \{(X_{i_j}, Y_{i_j}) : j = 1, \dots, M\}$$

and

$$F_n(U, Z) = D_n(U, Z) \setminus E_n(U, Z).$$

Then, given  $X_1, \dots, X_n$  and  $U$

$$F_n(U_+^k, Z) \equiv F_n(U_-^k, Z) \equiv F_n(U, Z),$$

while, because of

$$E_n(U_+^k, Z) = \{(X_{i_j}, W_{i_j}) : j = 1, \dots, M\}$$

and

$$E_n(U_-^k, Z) = \{(X_{i_j}, Z_{X_{i_j}}) : j = 1, \dots, M\},$$

under the event  $T_k$   $E_n(U_+^k, Z)$  and  $E_n(U_-^k, Z)$  have the same distribution. Since the event

$$\{(U_+^k, Z) \in A_n\} = \{\bar{R}_n(U_+^k, Z) \leq 1\}$$

also only depends on  $U \setminus U_k$  and  $Z_i$ 's falling out of block  $k$ ,

$$(F_n(U, Z), I_{\{(U_+^k, Z) \in A_n\}})$$

is independent from both  $E_n(U_+^k, Z)$  and  $E_n(U_-^k, Z)$ . This completes our statement.

We proceed with an indirect proof. Suppose that for all  $(u, z)$ ,  $\bar{R}_n(u, z) \rightarrow 0$ . This implies

$$\mathbf{P}\{(U, Z) \in A_n^c\} \rightarrow 0.$$

Thus, by Fatou's lemma

$$\begin{aligned} 0 &= \sup_{u, z} \limsup_n \frac{R_n(u, z)}{a_n} \geq \mathbf{E} \left\{ \limsup_n \frac{R_n(U, Z)}{a_n} \right\} \\ &\geq \mathbf{E} \{ \limsup_n a_n^{-1} R_n(U, Z) I_{\{(U, Z) \in A_n\}} \} \\ &\geq \limsup_n a_n^{-1} \mathbf{E} \{ R_n(U, Z) I_{\{(U, Z) \in A_n\}} \} \\ &= \limsup_n a_n^{-1} \mathbf{E} \{ |\phi_n(D_n(U, Z)) - L^*(U)| I_{\{(U, Z) \in A_n\}} \}. \end{aligned}$$

For proper  $k = k_n$ ,

$$\begin{aligned} &\limsup_n a_n^{-1} \mathbf{E} \{ |\phi_n(D_n(U, Z)) - L^*(U)| I_{\{(U, Z) \in A_n\}} \} \\ &\geq \limsup_n a_n^{-1} \mathbf{E} \left\{ \mathbf{E} \{ |\phi_n(D_n(U, Z)) - L^*(U)| \right. \\ &\quad \left. I_{\{(U, Z) \in A_n\}} | X_1, \dots, X_n, U \} I_{T_k} \right\} \\ &= \limsup_n a_n^{-1} \mathbf{E} \left\{ \left( \mathbf{E} \{ |\phi_n(D_n(U_-^k, Z)) - L^*(U_-^k)| \right. \right. \\ &\quad \left. \left. I_{\{(U_-^k, Z) \in A_n\}} | X_1, \dots, X_n, U \} I_{U_k=0} \right. \right. \\ &\quad \left. \left. + \mathbf{E} \{ |\phi_n(D_n(U_+^k, Z)) - L^*(U_+^k)| I_{\{(U_+^k, Z) \in A_n\}} | \right. \right. \\ &\quad \left. \left. X_1, \dots, X_n, U \} I_{U_k=1} \right) I_{T_k} \right\} \\ &= \limsup_n (2a_n)^{-1} \mathbf{E} \left\{ \left( \mathbf{E} \{ |\phi_n(D_n(U_-^k, Z)) - L^*(U_-^k)| \right. \right. \\ &\quad \left. \left. I_{\{(U_-^k, Z) \in A_n\}} | X_1, \dots, X_n, U \} \right. \right. \\ &\quad \left. \left. + \mathbf{E} \{ |\phi_n(D_n(U_+^k, Z)) - L^*(U_+^k)| I_{\{(U_+^k, Z) \in A_n\}} | X_1, \dots, X_n, U \} \right) I_{T_k} \right\} \\ &= \limsup_n (2a_n)^{-1} \mathbf{E} \left\{ \left( \mathbf{E} \{ |\phi_n(D_n(U_-^k, Z)) - L^*(U_-^k)| \right. \right. \\ &\quad \left. \left. I_{\{(U_-^k, Z) \in A_n\}} | X_1, \dots, X_n, U \} \right. \right. \\ &\quad \left. \left. + \mathbf{E} \{ |\phi_n(D_n(U_-^k, Z)) - L^*(U_-^k) - 1/2^{k+1}| \right. \right. \\ &\quad \left. \left. I_{\{(U_+^k, Z) \in A_n\}} | X_1, \dots, X_n, U \} \right) I_{T_k} \right\}, \end{aligned}$$

where in the last step we used that given  $X_1, \dots, X_n$  and  $U$  and the event  $T_k$ , the distributions of  $(D_n(U_-^k, Z), I_{\{(U_+^k, Z) \in A_n\}})$  and  $(D_n(U_+^k, Z), I_{\{(U_+^k, Z) \in A_n\}})$  are equal, and because of  $L^*(U_+^k) = L^*(U_-^k) + 1/2^{k+1}$ .

$$\begin{aligned} &\limsup_n a_n^{-1} \mathbf{E} \{ |\phi_n(D_n(U, Z)) - L^*(U)| I_{\{(U, Z) \in A_n\}} \} \\ &\geq \limsup_n (2a_n)^{-1} \mathbf{E} \left\{ \left( \mathbf{E} \{ |\phi_n(D_n(U_-^k, Z)) - L^*(U_-^k)| \right. \right. \\ &\quad \left. \left. I_{\{(U_-^k, Z) \in A_n\}} | X_1, \dots, X_n, U \} \right. \right. \\ &\quad \left. \left. + \mathbf{E} \{ |\phi_n(D_n(U_-^k, Z)) - L^*(U_-^k) - 1/2^{k+1}| \right. \right. \\ &\quad \left. \left. I_{\{(U_+^k, Z) \in A_n\}} | X_1, \dots, X_n, U \} \right) I_{T_k} \right\} \\ &\geq \limsup_n (2a_n)^{-1} \mathbf{E} \left\{ \mathbf{E} \{ 1/2^{k+1} \right. \\ &\quad \left. I_{\{(U_-^k, Z) \in A_n\}} I_{\{(U_+^k, Z) \in A_n\}} | X_1, \dots, X_n, U \} I_{T_k} \right\} \\ &= \limsup_n (2^{k+2} a_n)^{-1} \mathbf{P} \{ (U_-^k, Z) \in A_n, (U_+^k, Z) \in A_n, T_k \} \\ &= \limsup_n (2^{k+2} a_n)^{-1} \mathbf{P} \{ (U_-^k, Z) \in A_n, (U_+^k, Z) \in A_n \} \mathbf{P} \{ T_k \} \\ &\geq \limsup_n (2^{k+2} a_n)^{-1} \\ &\quad \left( 1 - 2 \frac{\mathbf{P} \{ (U_-^k, Z) \in A_n^c \} + \mathbf{P} \{ (U_+^k, Z) \in A_n^c \}}{2} \right) \mathbf{P} \{ T_k \} \\ &= \limsup_n (2^{k+2} a_n)^{-1} (1 - 2\mathbf{P} \{ (U, Z) \in A_n^c \}) \mathbf{P} \{ T_k \} \\ &\geq \limsup_n (2^{k+2} a_n)^{-1} (1 - 2\mathbf{P} \{ (U, Z) \in A_n^c \}) \left( 1 - \frac{n^2}{2^{2k+1} \ell_k} \right) \geq 1/2 \end{aligned}$$

if  $1/2^{k_n+2} \geq a_n$  (e.g., the choice  $k_n = \lfloor \log_2(1/a_n) \rfloor - 2 \rightarrow \infty$  is fine)

and  $\ell_{k_n} \geq n^2$ . This can be satisfied by the choice  $\ell_i = \max^2\{n : k_n \leq i\}$ .

This is a contradiction, therefore for all  $\{a_n\}$ , there is  $u, z, \{\ell_k\}$ , that  $\limsup_n \frac{R_n(u, z)}{a_n} > 0$ , which implies that for all  $\{a_n\}$ , there is  $u, z, \{\ell_k\}$ , that

$$\mathbf{E}|\phi_n(D_n(u, z)) - L^*(u)| > ca_n$$

infinitely often. Applying this to the original  $\sqrt{a_n}$ , this concludes the proof of the Theorem.  $\square$

### 3 CONCLUSION

In a standard pattern recognition design process, one takes a number of features, and evaluates whether these suffice or will do for discrimination. If so, a discrimination method is designed. If not, more or different features must be considered. The quality of a collection of features is measured by the Bayes probability of error,  $L^*$ . Thus, the first phase of any pattern recognition method is based on estimates of  $L^*$  (even *before* a discriminant is picked!). In this paper, we show that no one can trust any Bayes error estimate, and that it is futile to even let the sample size tend to infinity. It is only possible to give error bounds or confidence bands for  $L^*$  under assumptions on the distribution of the data. In practice, one can thus never claim to have a universally superior feature extraction or Bayes error estimation method, no matter how many simulations are performed and no matter how large the sample sizes are.

### ACKNOWLEDGMENTS

The second author's work was supported by NSERC Grant A3456 and by FCAR Grant 90-ER-0291. The first and the third author were supported by a grant from the Hungarian Academy of Sciences (MTA SZTAKI).

### REFERENCES

- [1] L. Birgé, "On Estimating a Density Using Hellinger Distance and Some Other Strange Facts," *Probability Theory and Related Fields*, vol. 71, pp. 271–291, 1986.
- [2] Z. Chen and K.S. Fu, "Nonparametric Bayes Risk Estimation for Pattern Classification," *Proc. IEEE Conf. Systems, Man, and Cybernetics*, Boston, 1973.
- [3] T.M. Cover, "Rates of Convergence for Nearest Neighbor Procedures," *Proc. Hawaii Int'l Conf. Systems Sciences*, pp. 413–415, Honolulu, 1968.
- [4] L. Devroye, "Any Discrimination Rule Can Have an Arbitrarily Bad Probability of Error for Finite Sample Size," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 4, pp. 154–157, 1982.
- [5] L. Devroye, "On Arbitrarily Slow Rates of Global Convergence in Density Estimation," *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, vol. 62, pp. 475–483, 1983.
- [6] L. Devroye, "Another Proof of a Slow Convergence Result of Birgé," *Statistics and Probability Letters*, vol. 23, pp. 63–67, 1995.
- [7] L. Devroye, L. Györfi, and G. Lugosi, *A Probabilistic Theory of Pattern Recognition*. New York/Berlin: Springer-Verlag, 1996.
- [8] K. Fukunaga and D.M. Hummels, "Bias of Nearest Neighbor Error Estimates," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 9, pp. 103–112, 1987.
- [9] K. Fukunaga and D.L. Kessel, "Estimation of Classification Error," *IEEE Trans. Computers*, vol. 20, pp. 1,521–1,527, 1971.
- [10] J.M. Garnett and S.S. Yau, "Nonparametric Estimation of the Bayes Error of Feature Extractors Using Ordered Nearest Neighbor Sets," *IEEE Trans. Computers*, vol. 26, pp. 46–54, 1977.
- [11] G. McLachlan, *Discriminant Analysis and Statistical Pattern Recognition*. New York: John Wiley, 1992.