

## ASYMPTOTICAL OPTIMALITY OF ADAPTIVE NEAREST NEIGHBOUR DISCRIMINATION\*

BY GUI-JING CHEN AND Y. H. WANG

*Anhui University, China and Concordia University, Canada*

This paper has discovered that the commonly used  $k - NN$  discrimination is not Bayesian risk consistent. A new adaptive discrimination procedure has been proposed. It has been proved that such procedure is Bayesian risk consistent, i.e. it has asymptotical optimality. Finally, an open problem is posed.

**1. Introduction, Procedure and Result.** Suppose that  $(X, \theta)$  is an  $R^d \times R^1$ -valued random variable, but  $\theta$  assumes only a finite number of values which, without loss of generality, can be denoted by  $1, 2, \dots, M$ . The intuitive background is that if an individual assumes  $j$  as its  $\theta$ -value, then we say that this individual belongs to class  $j$ . The problem is to determine the class  $j$  to which the individual is likely to belong, with the aid of its  $X$  value  $x$  and a series of i.i.d. observations  $Z^n = ((X_i, \theta_i), i = 1, 2, \dots, n)$  to be called the training sample.

Assume for the moment that the distribution of  $(X, \theta)$  is known. Then the conditional (priori) probabilities

$$P_j(x) = P(\theta = j \mid X = x), \quad j = 1, \dots, M. \quad (1)$$

can be calculated. Denote by  $\theta^*(x)$  the integer  $j^*$  such that

$$P_{j^*}(x) = \max \{P_j(x) : j = 1, \dots, M\}. \quad (2)$$

Then as known well,  $\theta^*$  attains the minimum error probability among all

---

\* Supported by National Natural Science Foundation of China and a grant from the Natural Science and Engineering Research Council of Canada.

AMS 1991 Subject Classifications: Primary 62H30, Secondary 62G20.

Key words and phrases: Adaptive procedure, Bayesian discrimination, Bayesian risk, nearest neighbour discrimination.

possible procedures, i.e.

$$R^* \triangleq P(\theta^*(X) \neq \theta) \tag{3}$$

$$= \inf\{P(\delta(X) \neq \theta) : \text{all } \delta\}. \tag{4}$$

The procedure  $\theta^*$  is called the Bayesian discrimination of the problem, its error probability  $R^*$  is called Bayesian risk.

Since the distribution of  $(X, \theta)$  is rarely known, a reasonable discrimination procedure must be based on the training sample  $Z^n$ . An intuitively attractive procedure, to be called the  $K - NN$  ( $K$ -Nearest Neighbour) discrimination, is defined as follows. Choose a distance  $\rho(\cdot, \cdot)$  on  $R^d$ , take a fixed integer  $k$  ( $1 \leq k < n$ ). For given  $x$ , reorder  $X_1, \dots, X_n$  according to the increasing order of  $\rho(x, X_i)$ ,  $i = 1, \dots, n$ :

$$\rho(x, X_{n1}) < \rho(x, X_{n2}) < \dots < \rho(x, X_{nk}) < \dots < \rho(x, X_{nn}),$$

where, if  $\rho(x, X_i)$  is the  $j$ -th smallest among  $\{\rho(x, X_t), t = 1, \dots, n\}$ , write  $(X_i, \theta_i) = (X_{nj}, \theta_{nj})$ , and call  $X_{nj}$  as  $j$ -th nearest neighbour point of  $x$  among  $X_t, t = 1, \dots, n$ . Think of  $\theta_{n1}, \dots, \theta_{nk}$  as  $k$  "ballots" for  $1, \dots, M$ . The "person" who receives the highest number of ballots is the choice (if more than one "person" receive the same highest number of ballots the choice may be decided randomly). This procedure, called  $k - NN$  discrimination, is denoted by  $\theta_n^{(k)}$  in the sequel,  $k$  is the order number of the discrimination.

$K - NN$  discrimination was first advanced by Fix and Hodges (1951), and has attracted much attention since the late sixties. Caver and Hart (1967), Wagner (1971), Fritz (1975), Devroye (1981), Chen, X.R. (1984), Bai, Z.D. (1985) and chen, G.J. and Kong, F.C. (1986) have studied the asymptotic behaviour of this procedure. The basic results are

$$\lim_{n \rightarrow \infty} R_n^{(k)} = R^{(k)}, \tag{5}$$

$$\lim_{n \rightarrow \infty} L_n^{(k)} = R^{(k)} \quad \text{a.s.}, \tag{6}$$

where  $R_n^{(k)}$  is error probability of the  $K - NN$  discrimination  $\theta_n^{(k)}$ :

$$R_n^{(k)} = P(\theta_n^{(k)}(X) \neq \theta); \tag{7}$$

and  $L_n^{(k)}$  is the empirical frequency of  $K - NN$  discrimination:

$$L_n^{(k)} = \frac{1}{n} \sum_{j=1}^n I(\theta^{(k)}(X_j) \neq \theta_j) \tag{8}$$

in which  $\theta^{(k)}$  is  $k - NN$  discrimination of  $\theta_j$  based on  $X_j$  as present sample, and  $(X_i, \theta_i)$ ,  $i \neq j$ ,  $i = 1, \dots, n$ , as training sample; and  $R^{(k)}$  is a positive number depending on  $k$  and discrimination of  $(X, \theta)$ .

We shall prove that (see Lemma 6 below) :  $R^* < R^{(k)}$  for every  $k$  fixed. It could be seen that  $k - NN$  discrimination is not Bayesian risk consistent, that is, it does not have asymptotical optimality. In order to improve this procedure, we consider an adaptive method. It is reasonable and natural to take

$$W_n^{(k)} = \sum_{j=1}^n (\theta_{j^n}^{(k)} \neq \theta_j) \tag{9}$$

as a basic cross-validatory criterion function to determine the order number of discrimination based on training sample  $Z^n$ . Define an adaptive order number as follows:

$$k^* = k^*(Z^n) = \max \left\{ k : W_n^{(k)} = \min (W_n^{(i)} : 1 \leq i \leq C_n) \right\}, \tag{10}$$

$$C_n = (\log n)^{\sigma_0}, \quad 0 < \sigma_0 < 1 \quad \text{fixed.} \tag{11}$$

For such order number  $k^*$ , we could operate  $k^* - NN$  procedure as above. About adaptive  $k^* - NN$  discrimination  $\theta_n^{(k^*)}$ , we have obtained some results as following.

**THEOREM.** *Under notations above, suppose that (i)  $X$  has a positive continuous density  $f(x)$ , for every sequence  $\varepsilon_n \downarrow 0$  there is a corresponding sequence  $b_n \uparrow \infty$  such that*

$$P(\|X\| > b_n) \leq \varepsilon_n \quad \text{and} \quad f(x) \geq \varepsilon_n^r \quad \text{as} \quad \|x\| \leq b_n \tag{12}$$

for some  $r > 0$ ;

(ii) *For  $j = 1, \dots, M$ , the conditional probabilities  $P_j(x)$  satisfy Lipschitz condition:*

$$\left| P_j(x) - P_j(x') \right| \leq D\rho(x, x'), \quad x, x' \in \mathbb{R}^d; \tag{13}$$

(iii) *There are  $i_0, j_0 (1 \leq i_0, j_0 \leq M)$  such that*

$$P(0 < P_{i_0}(X) < P_{j_0}(X)) > 0. \tag{14}$$

Then one has

$$\lim_{n \rightarrow \infty} P(\theta_n^{(k^*)} \neq \theta) = R^*, \tag{15}$$

$$\lim_{n \rightarrow \infty} L_n^{(k^*)} = R^* \quad \text{a.s.} \tag{16}$$

This theorem tells us that the adaptive  $K^* - NN$  discrimination is Bayesian risk consistent in both error sense.

**2. Proof of the Theorem.** The proof of the theorem consists of six lemmas. Let  $\xi_1, \dots, \xi_k$  be random variables such that  $\theta, \xi_1, \dots, \xi_k$  are i.i.d. under  $X = x$  given. Denote by  $\theta^{(k)}$  the discrimination value of  $\theta$  by "majority of votes" using  $\xi_1, \dots, \xi_k$ . Write

$$R^{(k)} = P(\theta^{(k)} \neq \theta). \tag{17}$$

**LEMMA 1.** One has

$$\lim_{k \rightarrow \infty} R^{(k)} = R^*. \tag{18}$$

**PROOF.** From (2), (3)

$$R^* = 1 - P(\theta^*(X) = \theta) = 1 - EP_{j^*}(X) \tag{19}$$

and from (17)

$$\begin{aligned} R^{(k)} &= 1 - \sum_{u=1}^M P(\theta^{(k)} = u, \theta = u) \\ &= 1 - \sum_{u=1}^M E\left\{P(\theta^{(k)} = u | X)P(\theta = u | X)\right\}. \end{aligned} \tag{20}$$

Write  $N_u^{(k)} = \sum_{i=1}^k I(\xi_i = u)$ ,  $u = 1, \dots, M$ . One has

$$\lim_{k \rightarrow \infty} P\left\{\sqrt{k}\left[\left(N_u^{(k)} - N_{j^*}^{(k)}\right)/k - \left(P_u(x) - P_{j^*}(x)\right)\right]/\sigma \leq t | x\right\} = \Phi(t)$$

by central limit theorem, where  $\sigma^2 = P_u(x) + P_{j^*}(x) - (P_u(x) + P_{j^*}(x))^2$ . Then, if  $P_{j^*}(x) > P_u(x)$ ,

$$0 \leq P(P = u | x) \leq P(N_u^{(k)} \geq N_{j^*}^{(k)} | x) \rightarrow 0 \quad \text{as } k \rightarrow \infty$$

by the definition of  $\theta^{(k)}$ . From (19) and (20),

$$\begin{aligned} \lim_{k \rightarrow \infty} P(\theta^{(k)} = j^* | x) &= 1, \\ \lim_{k \rightarrow \infty} R^{(k)} &= 1 - \sum_{u=1}^M E\left\{\lim_{k \rightarrow \infty} P(\theta^{(k)} = u | X)P(\theta = u | X)\right\} = R^*. \end{aligned}$$

LEMMA 2. Under the assumptions (i) and (ii), one has

$$\lim_{n \rightarrow \infty} R_n^{(k)} = R^{(k)} \tag{21}$$

uniformly for  $k \leq (\log n)^{\sigma_0}$ .

PROOF.

$$\begin{aligned} R_n^{(k)} &= 1 - \sum_{u=1}^M E \left\{ P(\theta_n^{(k)} = u, \theta = u \mid X, X_1, \dots, X_n) \right\} \\ &= 1 - \sum_{u=1}^M E \left\{ P(\theta_n^{(k)} = u \mid X, X_{n1}, \dots, X_{nk}) P(\theta = u \mid X) \right\}. \end{aligned}$$

From this and (20)

$$\left| R_n^{(k)} - R^{(k)} \right| \leq \sum_{u=1}^M E \left| P(\theta_n^{(k)} = u \mid X, X_{n1}, \dots, X_{nk}) - P(\theta^{(k)} = u \mid X) \right|. \tag{22}$$

By the definition of  $\theta_n^{(k)}$ , one has

$$\begin{aligned} &P(\theta_n^{(k)} = u \mid x, X_{n1}, \dots, X_{nk}) \\ &= \sum^* \prod_{m=1}^K P_u(X_{nm})^{\Delta_{mn}} \prod_{t=1}^{M-1} P_{i_t}(X_{um})^{\Delta_{mi_t}}, \end{aligned} \tag{23}$$

where  $\sum^* = \sum^{(1)} \frac{1}{v} \sum^{(2)} \sum^{(3)} \sum^{(4)} \sum^{(5)}$ ,

$$\begin{aligned} \sum^{(1)} &= \sum_{v=1}^k, & \sum^{(2)} &= \sum_{s=1}^{k/v}, \\ \sum^{(3)} &= \sum (i_1, \dots, i_{M-1}) \in A_u, \\ \sum^{(4)} &= \sum (K_l, K_{i_1}, \dots, K_{i_{M-1}}) \in B, \\ \sum^{(5)} &= \sum (\Delta_{mu}, \Delta_{mi_t} : m = 1, \dots, k; t = 1, \dots, M-1) \in C \end{aligned}$$

in which

$$\begin{aligned} A_u &= \{(i_1, \dots, i_{M-1}) : \text{all permutations of } (1, \dots, u-1, u+1, \dots, M)\}, \\ B &= B(i_1, \dots, i_{M-1}) \\ &= \left\{ (K_u, K_{i_1}, \dots, K_{i_{M-1}}) : \sum_{i=1}^M K_i = K, K_u = K_{i_t} = s, t = 1, \dots, v-1; \right. \end{aligned}$$

$$\begin{aligned}
& \left. k_{i_t} < s, t = v, v+1, \dots, M-1 \right\}, \\
C &= C(k_u, k_{i_1}, \dots, k_{i_{M-1}}) \\
&= \left\{ (\Delta_{m_u}, \Delta_{m_{i_t}} : m = 1, \dots, k, t = 1, \dots, M-1) : \Delta = 0, \text{ or } 1, \right. \\
& \quad \left. \sum_{i=1}^M \Delta_{m_i} = 1, \sum_{n=1}^k \Delta_{m_i} = k_i \right\}.
\end{aligned}$$

One could also get

$$P(\theta^{(k)} = u | x) = \sum^* \prod_{m=1}^k (Pu(x))^{\Delta_{m_u}} \prod_{t=1}^{M-1} P_{i_t}(x)^{\Delta_{m_{i_t}}}. \quad (24)$$

Then we obtain

$$\begin{aligned}
& \left| P(\theta_n^{(k)} = u | x, X_{n1}, \dots, X_{nk}) - P(\theta^{(k)} = u | x) \right| \\
& \leq \sum^* \sum_{m=1}^K \sum_{i=1}^M \left| P_i(X_{nm}) - P_i(x) \right|. \quad (25)
\end{aligned}$$

Using the assumptions (i) and (ii), one could prove that

$$\begin{aligned}
& P\left( \left| P_i(X_{nm}) - P_i(x) \right| \geq \varepsilon/N | x \right) \leq (k+1)n^k \exp(-cf(x)\varepsilon^d n^{1-\sigma}), \\
& P\left\{ \left| P(\theta_n^{(k)} = u | x, X_{n1}, \dots, X_{nk}) - P(\theta^{(k)} = u | x) \right| \geq \varepsilon | x \right\} \\
& \leq N(k+1)n^k \exp(-cf(x)\varepsilon^d n^{1-\sigma}). \quad (26)
\end{aligned}$$

where  $N = M^{ck}$ ,  $C$ ,  $0 < \sigma < 1$  stand for positive constant independent of  $n$ ,  $k \leq (\log n)^{\sigma_0}$ .

Take  $G > 0$  large enough such that  $P(\|X\| > G) < \varepsilon$ . Define  $\delta = \inf \left\{ f(x) : \|x\| \leq G \right\} > 0$ . Then for  $\|x\| \leq G$

$$\begin{aligned}
& P\left( \left| P(\theta_n^{(k)} = u | x, X_{n1}, \dots, X_{nk}) - P(\theta^{(k)} = u | x) \right| \geq \varepsilon | x \right) \\
& \leq N(k+1)n^k \exp(-C\delta\varepsilon^d n^{1-\sigma}) \leq C \exp(-Cn^{1-\sigma})
\end{aligned}$$

when  $n$  large enough,  $k \leq (\log n)^{\sigma_0}$ . Then

$$\begin{aligned}
& E \left| P(\theta_n^{(k)} = u | X, X_{n1}, \dots, X_{nk}) - P(\theta^{(k)} = u | X) \right| \\
& \leq \varepsilon + E \left\{ \left| P(\theta_n^{(k)} = u | X, X_{n1}, \dots, X_{nk}) - P(\theta^{(k)} = u | X) \right| I(\|X\| \leq G) \right\} \\
& \leq \varepsilon + \varepsilon + E \left\{ P\left( \left| P(\theta_n^{(k)} = u | X, X_{n1}, \dots, X_{nk}) \right. \right. \right. \\
& \quad \left. \left. \left. - P(\theta^{(k)} = u | X) \right| \geq \varepsilon | X \right) I(\|X\| \leq G) \right\} \\
& \leq 2\varepsilon + c \exp(-Cn^{1-\sigma}) < 3\varepsilon
\end{aligned}$$

as  $n$  large enough for  $k \leq (\log n)^{\sigma_0}$ . From this and (22), Lemma 2 is proved.

**LEMMA 3.** *Under the assumptions (i) and (ii), for any bound measurable function  $\mu(x)$ ,  $x \in R^d$ , one has*

$$P\left\{\frac{1}{n} \sum_{j=1}^n \left| \mu(X_{j^{n_i}}) - \mu(X_j) \right| \geq \varepsilon_n\right\} \leq C \exp(-bn^{1-\sigma}), \quad (27)$$

$i = 1, \dots, k, k \leq (\log n)^{\sigma_0}$ ,  $n$  large enough. Where  $X_{j^{n_i}}$  is  $i$ -th nearest neighbour point of  $X_j$ , among  $X_t : t \neq j, t = 1, \dots, n; \varepsilon_n$  satisfies

$$\varepsilon_n \downarrow 0, \quad \varepsilon_n N^r \rightarrow \infty \quad \text{as } n \rightarrow \infty \quad \text{for any } r > 0; \quad (28)$$

and  $c, b, \sigma \in (0, 1)$  are independent of  $n$ .

**PROOF.** At first suppose that  $\mu(x) = I_A(x)$ , where  $A$  is a rectangle in  $R^d$ . Denote the boundary of  $A$  by  $\partial A$ , and denote the probability distribution or measure of  $X$  by  $F$ . Because  $F(\partial A) = 0$ , so we could choose rectangles  $A_{1n}, A_{2n}$  such that

$$A_{1n} \subset A \subset A_{2n},$$

$$0 < \rho(\partial A_{1n}, \partial A_{2n}) \triangleq \rho_n \rightarrow 0 \quad \text{as } n \rightarrow \infty, \quad (29)$$

$$F(A_{2n} \cap A_{1n}^c) \triangleq \delta_n \quad (30)$$

in which  $\rho(\partial A_{1n}, \partial A_{2n})$  is distance between  $\partial A_{1n}$  and  $\partial A_{2n}$ . Choose  $\delta_n$  such that

$$\varepsilon_n/8 \leq \delta_n \leq \varepsilon_n/4. \quad (31)$$

By the condition (12), there are rectangles  $A_{3n}$  such that  $A_{2n} \subset A_{3n}$  and

$$F(A_{3n}^c) \leq \varepsilon_n/4, \quad f_n \triangleq \inf \{f(x), x \in A_{3n}\} \geq \varepsilon_n^r. \quad (32)$$

Then the set  $A_{1n} \cup (A_{2n}^c \cap A_{3n})$  could be divided into  $N_1$  rectangles  $B_1, \dots, B_{N_1}$  such that  $B_i \cap B_j = \emptyset, i \neq j$ , and  $N_1 \leq C\rho_n^{-d}$ ,

$$d(B_i) \leq \rho_n/2, \quad F(B_i) \geq cf_n\rho_n^d, \quad i = 1, \dots, N_1, \quad (33)$$

here  $d(B_i)$  is diameter of  $B_i$ . Then

$$P_n \triangleq \inf \{F(B_i), i = 1, \dots, N_1\} \geq c\varepsilon_n^{r'} \quad (r' = r + d) \quad (34)$$

by (31), (32). Denote

$$S_1 = \bigcap_{i=1}^{N_1} \left\{ \sum_{j=1}^n I(X_j \in B_i) \geq k \right\};$$

$$S_2 = \left\{ \sum_{j=1}^n I(X_j \in (A_{1n}^c \cap A_{2n}) \cup A_{3n}^c) \leq n\varepsilon_n \right\}.$$

From (30), (31) and (32), one has

$$P\{A_{1n}^c \cap A_{2n} \cup A_{3n}^c\} \leq \varepsilon_n/2.$$

Using Bennett inequality (Bennett, G. (1962)) and condition (28)

$$P(S_2^c) \leq 2 \exp(-2n\varepsilon_n^2) \leq c \exp(-bn^{1-\sigma}). \quad (35)$$

Because  $k \leq (\log n)^{\sigma_0}$ ,  $N_1 \leq C\varepsilon_n^{-d}$ , (28) and (34), one has

$$P(S_1^c) \leq \sum_{i=1}^{N_1} \sum_{j=0}^{k-1} C_n^j (F(B_i))^j (1 - F(B_i))^{n-j} \leq c \exp(-bn^{1-\sigma}). \quad (36)$$

By the definition of  $S_1$ , for  $j = 1, \dots, n$ ;  $i = 1, \dots, k$ , one has

$$S_1 \cap \{X_j \in A\} \Rightarrow \{X_{jni} \in A\},$$

$$S_1 \cap \{X_j \in A_{3n} \cap A_{2n}^c\} \Rightarrow \{X_{jni} \notin A\}.$$

Therefore when the event  $S_1$  happens, then

$$|I_A(X_j) - I_A(X_{jni})| \leq I(X_j \in (A_{2n} \cap A_{1n}^c) \cup A_{3n}^c).$$

Then one has

$$P\left\{ \frac{1}{n} \sum_{j=1}^n |I_A(X_j) - I_A(X_{jni})| \geq \varepsilon_n \right\}$$

$$\leq P\left\{ \frac{1}{n} \sum_{j=1}^n I(X_j \in A_{3n}^c \cup (A_{2n} \cap A_{1n}^c)) > \varepsilon_n \right\} + P\{S_1^c\}$$

$$\leq c \exp(-bn^{1-\sigma})$$

by (35) and (36). Then (27) holds when  $\mu(x) = I_A(x)$ . Using commonly used measure theory method, it could be shown further that (27) is still true for a general bound measurable function  $\mu(x)$ . Lemma 3 is proved.



LEMMA 4. Suppose that  $X_{jn_1}, \dots, X_{jn_k}$  are the first  $k$  nearest neighbour points of  $X_j$ , among  $X_t : t \neq j, t = 1, \dots, n$ . Then with probability one, we could divide the set class  $J^{(n)} = \{(j, j_{n_1}, \dots, j_{n_k}) : j = 1, \dots, n\}$  into  $q$  (independent of  $n$ ) subclasses  $J_1^{(n)}, \dots, J_q^{(n)}$  such that for each  $J_t^{(n)}$ , when  $j' \neq j$  and  $(j, j_{n_1}, \dots, j_{n_k}) \in J_t^{(n)}, (j', j'_{n_1}, \dots, j'_{n_k}) \in J_t^{(n)}$

$$\{j, j_{n_1}, \dots, j_{n_k}\} \cap \{j', j'_{n_1}, \dots, j'_{n_k}\} = \emptyset. \tag{37}$$

(See the Lemma 2 of [9].)

LEMMA 5. Under the assumptions (i) and (ii), one has

$$P(|L_n^{(k)} - R^{(k)}| \geq \varepsilon_n) \leq c \exp(-bn^{1-\sigma}) \tag{38}$$

for  $k \leq (\log n)^{\sigma_0}$ ,  $n$  large enough. Where  $L_n^{(k)}, R^{(k)}$  are defined by (8) and (17);  $\{\varepsilon_n\}$  satisfies the condition (28);  $0 < c, b, 0 < \sigma < 1$  are independent of  $n$ .

PROOF. Let  $\xi_i, i = 1, \dots, k, \theta^{(k)}$  be defined as in the Lemma 1. Then by (8) and (17),

$$|L_n^{(k)} - R^{(k)}| \leq \sum_{u=1}^M \Delta_{nk}^{(u)},$$

where

$$\begin{aligned} \Delta_{nk}^{(u)} &\leq \left| \frac{1}{n} \sum_{j=1}^n I(\theta_{jn}^{(k)} = u, \theta_j = u) - T_{nk1}^{(u)} \right| \\ &\quad + \left| T_{nk1}^{(u)} - T_{nk2}^{(u)} \right| + \left| T_{nk2}^{(u)} - P(\theta^{(k)} = u, \theta = u) \right| \\ &\triangleq \Delta_{uk1}^{(u)} + \Delta_{nk2}^{(u)} + \Delta_{nk3}^{(u)}, \end{aligned} \tag{39}$$

where

$$T_{nk1}^{(u)} = \frac{1}{n} \sum_{j=1}^n P^{(n,k,u)}(X_{jn_1}, \dots, X_{jn_k}) P_u(X_j) \tag{40}$$

in which

$$\begin{aligned} P^{(n,k,u)}(X_{jn_1}, \dots, X_{jn_k}) &= P(\theta_{jn}^{(k)} = u \mid X_{jn_1}, \dots, X_{jn_k}), \\ T_{nk2}^{(u)} &= \frac{1}{n} \sum_{j=1}^n P^{(n,k,u)}(X_j, \dots, X_j) P_u(X_j). \end{aligned}$$

Let

$$\begin{aligned} \mu(X) &= P^{(n,k,u)}(X, \dots, X) P_u(X), \tag{41} \\ N_{jni}^{(k)} &= \sum_{t=1}^k I(\theta_{jnt} = i), \quad i = 1, \dots, M, \end{aligned}$$

Where  $\theta_{jnt}$  is matched with  $X_{jnt}$ , which is  $t$ -th nearest neighbour point of  $X_j$ , among  $X_s : s \neq j, s = 1, \dots, n$ . Because

$$P(\theta_{jnt} = i \mid X_{jn1}, \dots, X_{jnk}) = P_i(X_{jnt}),$$

then

$$P^{(n,k,u)}(X_{jn1}, \dots, X_{jnk}) = \sum^* \prod_{m=1}^k P_u(X_{jnm})^{\Delta_{mu}} \prod_{t=1}^{M-1} P_{i_t}(X_{jnm})^{\Delta_{mi_t}}, \tag{42}$$

where summation  $\sum^*$  is defined as (23). So

$$\begin{aligned} & |P^{(n,k,u)}(X_{jn1}, \dots, X_{jnk}) - P^{(n,k,u)}(X_j, \dots, X_j)| \\ & \leq \sum^* \sum_{m=1}^k \sum_{i=1}^M |P_i(X_{jnm}) - P_i(X_j)|. \end{aligned} \tag{43}$$

Let  $N$  be defined as Lemma,  $\varepsilon'_n = \varepsilon_n/N$ . For  $k \leq (\log n)^{\sigma_0}$ ,  $\{\varepsilon'_n\}$  still satisfies (28). Using Lemma 3, one has

$$\begin{aligned} P(\Delta_{nk2}^{(u)} > \varepsilon_n) & \leq \sum^* \sum_{m=1}^k \sum_{i=1}^M P\left(\frac{1}{n} \sum_{j=1}^n |P_i(X_{jnm}) - P_i(X_j)| > \varepsilon'_n\right) \\ & \leq N \exp(-bn^{1-\sigma}) \leq c \exp(-bn^{1-\sigma}). \end{aligned} \tag{44}$$

From (24), one has

$$\begin{aligned} P(\theta^{(k)} = u \mid X) & = \sum^* \prod_{m=1}^k P_u(X)^{\Delta_{mu}} \prod_{t=1}^{M-1} P_{i_t}(X)^{\Delta_{mi_t}} \\ & = P^{(n,k,u)}(X, \dots, X). \end{aligned} \tag{45}$$

From this and (41),  $\mu(x) = P(\theta^{(k)} = u \mid x)P(\theta = u \mid x)$ , which is independent of  $n$ , and

$$\begin{aligned} P(\theta^{(k)} = u, \theta = u) & = E\{P(\theta^{(k)} = u, \theta = u \mid X)\} \\ & = E\{P(\theta^{(k)} = u \mid X)P(\theta = u \mid X)\} = E\mu(X). \end{aligned}$$

From (39),

$$\Delta_{nk3}^{(u)} = \left| \frac{1}{n} \sum_{j=1}^n \mu(X_j) - E\mu(X) \right|,$$

using Bennett inequality (Bennett, G. (1962)),

$$P(\Delta_{nk3}^{(u)} > \varepsilon_n) \leq 2 \exp(-bn^{1-\sigma}). \tag{46}$$

Now consider the term  $\Delta_{nk1}^{(u)}$ . By Lemma 4, one could divide the set class  $J^{(n)} = \{(j, j_{n1}, \dots, j_{nk}) : j = 1, \dots, n\}$  into  $q$  subclasses  $J_t^{(n)}$ ,  $t = 1, \dots, q$  such that (37). Denote

$$\cup(j, j_{n1}, \dots, j_{nk}) = I(\theta_{jn}^{(k)} = u, \theta_j = u) - P^{(n,k,u)}(X_{jn1}, \dots, X_{jnk}).$$

Then we have

$$\begin{aligned} \Delta_{nk1}^{(u)} &= \left| \frac{1}{n} \sum_{j=1}^n \cup(j, j_{n1}, \dots, j_{nk}) \right| \\ &= \left| \sum_{t=1}^q \frac{1}{n} \sum_{(j, j_{n1}, \dots, j_{nk}) \in J_t^{(n)}} \cup(j, j_{n1}, \dots, j_{nk}) \right|. \end{aligned}$$

Using Bennett inequality (Bennett, G. (1962)) again

$$P(\Delta_{nk1}^{(u)} > \varepsilon_n \mid X_1, \dots, X_n) \leq 2q \exp(-2n\varepsilon_n^2/q^2) \leq c \exp(-bn^{1-\sigma}).$$

Then

$$P(\Delta_{nk1}^{(u)} > \varepsilon_n) = EP(\Delta_{nk1}^{(u)} > \varepsilon_n \mid X_1, \dots, X_n) \leq c \exp(-bn^{1-\sigma}). \tag{47}$$

Combining (39), (44), (46) and (47), this lemma is proved.

**LEMMA 6.** Under the assumptions (i), (ii) and (iii), one has

$$R^* < R^{(k)}, \quad k = 1, 2, \dots \tag{48}$$

$$\lim_{n \rightarrow \infty} K^* = \infty \quad \text{a.s.}, \tag{49}$$

where  $R^*, R^{(k)}$  and  $K^*$  are defined by (3), (10) and (17).

**PROOF.** Write  $A = \{x : 0 < P_{i_0}(x) < P_{j_0}(x)\}$ . From (2),  $P_i(x) \leq P_{j^*}(x)$  for all  $x \in R^d$  and all  $i = 1, \dots, M$ . When  $x \in A$ , then  $0 < P_{i_0}(x) < P_{j^*}(x)$ ,  $P(\theta^{(k)} = i_0 \mid x) \geq (P_{i_0}(x))^k > 0$  by the definition of  $\theta^{(k)}$ , therefore

$$\sum_{i=1}^M P(\theta^{(k)} = i \mid x) P_i(x) < \sum_{i=1}^M P(\theta^{(k)} = i \mid x) P_{j^*}(x) = P_{j^*}(x).$$

Because  $P(X \in A) > 0$ , then

$$R^* = 1 - EP_{j^*}(X) < 1 - \sum_{i=1}^M EP(\theta^{(k)} = i \mid X) P_i(X) = R^{(k)}$$

by (19) and (20). So (54) is true.

Denote

$$\begin{aligned} k_n &= [(\log n)^{\sigma_0}], \\ B_0 &= \{L_n^{(k_n)} - R^{(k_n)} \rightarrow 0\}, \\ B_k &= \{L_n^{(k)} - R^{(k)} \rightarrow 0\}, \quad k = 1, 2, \dots, \\ B &= \bigcap_{K=0}^{\infty} B_k. \end{aligned}$$

Using (38) one has

$$\sum_{n=1}^{\infty} P(|L_n^{(k)} - R^{(k)}| \geq \varepsilon) < \infty, \quad \sum_{n=1}^{\infty} P(|L_n^{(k_n)} - R^{(k_n)}| > \varepsilon) < \infty,$$

then  $P(B_k) = 1$ ,  $k = 0, 1, 2, \dots$ , by Borel-Cantelli lemma. So  $P(B) = 1$ . We come to prove that for  $((X_i, \theta_i), i = 1, \dots, n, \dots) \in B$ , one has

$$K^*(Z^n) \rightarrow \infty, \quad \text{as } n \rightarrow \infty. \tag{50}$$

In fact, if it is not true, then there is a subsequence  $\{n_i\}$  such that as  $i \rightarrow \infty$ ,  $n_i \rightarrow \infty$ , but

$$K^*(Z^{n_i}) \leq k_0 < \infty. \tag{51}$$

Because  $B \subset B_k$ , one has

$$\liminf_{i \rightarrow \infty} L_{n_i}^{(k^*(Z^{n_i}))} \geq \min_{1 \leq k \leq k_0} R^{(k)} \tag{52}$$

by the definition of  $B_k$  and (51). On other hand,  $B \subset B_0$ ,

$$\limsup_{i \rightarrow \infty} L_{n_i}^{(k^*(Z^{n_i}))} \leq \limsup_{i \rightarrow \infty} L_{n_i}^{(k_{n_i})} = R^* \tag{53}$$

by (10), Lemma 1 and 5. It is easy to see that (52), (53) and (48) are contradictory. Lemma 6 is proved.

Now we could give the proof of the theorem as follows. For  $\varepsilon > 0$  given arbitrarily, there is  $k_0$  large enough such that when  $k \geq k_0$  then  $|R^{(k)} - R^*| < \varepsilon/3$  by Lemma 1; and there is  $n_1 (\geq K_0)$  large enough such that when  $n \geq n_1$  then  $P(k^* \leq k_0) < \varepsilon/3$  by Lemma 6; and by Lemma 2 there is  $n_2 (\geq n_1)$  large enough such that as  $n \geq n_2$  then  $|R_n^{(k)} - R^{(k)}| < \varepsilon/3$ , for  $k \leq (\log n)^{\sigma_0}$ .

Therefore as  $n \geq n_2$ , one has

$$\begin{aligned} |P(\theta_n^{(k^*)} \neq \theta) - R^*| &\leq \sum_{1 \leq k \leq (\log n)^{\sigma_0}} |R^{(k)} - R^*| P(k^* = k) \\ &\leq P(k^* \leq k_0) + \sum_{K_0 \leq K \leq (\log n)^{\sigma_0}} (|R_n^{(k)} - R^{(k)}| + |R^{(k)} - R^*|) P(k^* = k) \\ &< \varepsilon/3 + \varepsilon/3 + \varepsilon/3 = \varepsilon \end{aligned}$$

hence (15) is true. In order to prove (16), denote  $\varepsilon_n = \varepsilon/(\log n)^{\sigma_0}$ , which satisfies the condition (28). From Lemma 5,

$$\begin{aligned} P(|L_n^{(k^*)} - R^{(k^*)}| > \varepsilon) &\leq \sum_{1 \leq K \leq (\log n)^\sigma} P(|L_n^{(k)} - R^{(k)}| > \varepsilon_n) \\ &\leq \sum_{1 \leq K \leq (\log n)^\sigma} P(|L_n^{(k)} - R^{(k)}| > \varepsilon_n) \\ &\leq (\log n)^{\sigma_0} c \exp(-bn^{1-\sigma}) \leq c \exp(-bn^{1-\sigma}) \end{aligned}$$

Then

$$\sum_{n=1}^{\infty} P(|L_n^{(k^*)} - R^{(k^*)}| > \varepsilon) < \infty,$$

therefore

$$L_n^{(k^*)} - R^{(k^*)} \rightarrow 0, \quad \text{as } n \rightarrow \infty \quad \text{a.s.}$$

By Borel Cantelli lemma. Using Lemmas 1 and 6, one has

$$R^{(k^*)} - R^* \rightarrow 0, \quad \text{a.s. as } n \rightarrow \infty,$$

it implies (16). The theorem is proved completely.

**3. Discussion.** From intuition point of view, in order to improve the convergence rate we should take  $\sigma_0$  as large as possible when choosing adaptive discrimination order number  $k^*$  by (10) and (11). In fact it is not difficult to justify that if we take  $\sigma_n = 1 - (\log \log n)^{-\sigma_1}$  ( $0 < \sigma_1 < 1$ ), instead of  $\sigma_0$ , the theorem still holds. But it seems to us that  $\sigma_0$  could not be greater than or equal to 1. We now do not know whether it is true.

## REFERENCES

- FIX, E. and HODGES, J. (1951). Randolph Field, Texas, Project 21-49-004.
- COVER, T. M. and HART, P. E. (1967). Nearest neighbour pattern classification, *IEEE, Trans. Inform. Theory* 21-27.
- WAGNER, T. J. (1971). Convergence of the nearest neighbour rule, *IEEE, Trans. Inform. Theory* 1971, 566-570.
- FRITZ, J. (1975). Distribution-free exponential error bound for nearest neighbour pattern classification, *IEEE, Trans. Inform. Theory* 552-557.
- DEVROYE, L. (1981). On the almost everywhere convergence of nonparametric regression function estimates, *Ann. Statist.* 1320-1327.

CHEN, X. R. (1984). Exponential posterior error bound for the  $k - NN$  discrimination rule, *Scientia Sinica* 11(A), 978–986.

BAI, Z. D. (1985). The strong consistency of error probability estimates in  $NN$  discrimination, *Chinese Ann. Math.* 299–308.

CHEN, G. J. and KONG, F. C. (1986). Sufficient and necessary condition for convergence of conditional error probability in  $NN$  discrimination, *J. Math. Research and Exposition* 2, 95–104.

CHEN, G. J. (1986). On estimation of conditional risk of  $k - NN$  prediction, *J. of Anhui Univ.* 1, 1–9.

BENNETT, G. (1962). Upper bounds on probability inequalities for the sum of independent, bounded random variables, *J. Amer. Statist. Assoc.* 33–45.

DEPARTMENT OF MATHEMATICS  
CONCORDIA UNIVERSITY  
MONTREAL, H3G 1M8  
QUEBEC, CANADA

DEPARTMENT OF MATHEMATICS  
ANHUI UNIVERSITY  
HEFEI 230039  
ANHUI, CHINA