# ON THE STRONG UNIVERSAL CONSISTENCY OF NEAREST NEIGHBOR REGRESSION FUNCTION ESTIMATES[1]

By Luc Devroye, László Györfi, Adam Krzyżak and Gábor Lugosi

*McGill University, Technical University of Budapest, Concordia University and Technical University of Budapest*

Two results are presented concerning the consistency of the $k$–nearest neighbor regression estimate. We show that all modes of convergence in $L_1$ (in probability, almost sure, complete) are equivalent if the regression variable is bounded. Under the additional condition $k/\log n \to \infty$ we also obtain the strong universal consistency of the estimate.

**1. Introduction.** Let $(X_1, Y_1), \ldots, (X_n, Y_n)$ be independent observations of an $\mathcal{R}^d \times \mathcal{R}$-valued random vector $(X, Y)$. Denote the probability measure of $X$ by $\mu$. The regression function $m(x) = \mathbf{E}(Y \mid X = x)$ can be estimated by the *kernel estimate*

$$m_n(x) = \frac{\sum_{i=1}^n Y_i K_h(x - X_i)}{\sum_{i=1}^n K_h(x - X_i)},$$

where $h > 0$ is a smoothing factor depending upon $n$; $K$ is an absolutely integrable function (the kernel); and $K_h(x) = K(x/h)$ [Nadaraya (1964, 1970), Watson (1964)]. Alternatively, one can use the $k$–nearest neighbor estimate,

$$m_n(x) = \sum_{i=1}^n W_{ni}(x; X_1, \ldots, X_n) Y_i,$$

and $W_{ni}(x; X_1, \ldots, X_n)$ is $1/k$ if $X_i$ is one of the $k$ nearest neighbors of $x$ among $X_1, \ldots, X_n$, and $W_{ni}$ is zero otherwise. Note in particular that $\sum_{i=1}^n W_{ni} = 1$. The $k$–nearest neighbor estimate was studied by Cover (1968). For a survey of other estimates, see, for example, Collomb (1981, 1985) or Györfi (1981).

We are concerned with the $L_1$ convergence of $m_n$ to $m$ as measured by $J_n = \int |m_n(x) - m(x)| \mu(dx)$, where $\mu$ is the (unknown) probability measure for $X$. This quantity is particularly important in discrimination based on the kernel rule [see Devroye and Wagner (1980) or Stone (1977)]. Stone (1977) first pointed out that there exist estimators for which $J_n \to 0$ in probability for all distributions of $(X, Y)$ with $\mathbf{E}|Y| < \infty$. This included the nearest neighbor and histogram estimates. For example, for the $k$ nearest neighbors, it suffices to ask that

(1) $$k \to \infty, \qquad k/n \to 0,$$

provided that ties among points at equal distance from $x$ are adequately taken care of. These conditions are the best possible. Devroye and Wagner (1980) and, independently, Spiegelman and Sacks (1980) showed that this is also the case for the kernel estimate with smoothing factor $h$ provided that $K$ is a bounded nonnegative function with compact support such that, for a small fixed sphere $S$ centered at the origin, $\inf_{x \in S} K(x) > 0$, and that

$$(2) \qquad\qquad \lim_{n \to \infty} h = 0, \qquad \lim_{n \to \infty} nh^d = \infty.$$

These results were extended and complemented by Greblicki, Krzyżak and Pawlak (1984), Krzyżak (1986) and Krzyżak and Pawlak (1984).

Interestingly, it turns out that the conditions for the "in probability" convergence of $J_n$ are also sufficient for the strong convergence of $J_n$, thus rendering all modes of convergence equivalent. Difficulties arise when the $X$-variable does not have an absolutely continuous distribution. We summarize what is known in this respect:

1. For the $k$–nearest neighbor estimates, $J_n \to 0$ almost surely under condition (1) whenever $X$ has a density and $Y$ is bounded [Devroye and Györfi (1985), Chapter 10, and Zhao (1987)]. Beck (1979) showed this result earlier under the additional constraint that $m$ has a continuous version.
2. For the $k$–nearest neighbor estimate, $J_n \to 0$ almost surely for all distributions of $(X, Y)$ with $Y$ bounded, provided that $k/n \to 0$ and $k/\log \log n \to \infty$ [Devroye (1982)]. The unnatural condition on $k$ arises from the proof method: the convergence of $J_n$ to 0 is obtained by first establishing the pointwise convergence (i.e., $m_n - m \to 0$ almost surely) at almost all $x(\mu)$ and then moving on to $L_1$ convergence via a result of Glick (1974).
3. Devroye and Györfi (1983) obtained the equivalence for all distributions of $(X, Y)$ with $|Y| \le M < \infty$ for the *histogram regression estimate*. Györfi (1991) has pointed out that $J_n \to 0$ almost surely for a modification of *partitioning estimates* whenever $\mathbf{E}|Y| < \infty$, provided that a bin width condition similar to (2) [with the additional condition $nh^d/\log(n) \to \infty$] is satisfied.
4. Assuming that $Y$ is uniformly bounded, the *kernel estimate* is strongly consistent if (2) holds, $K$ is a Riemann integrable kernel and $K \ge aI_S$, where $a > 0$ is a constant and $S$ is a ball centered at the origin that has a positive radius [Devroye and Krzyżak (1989)].

The purpose of the present paper is twofold. First we explain a simple technique based upon exponential martingale inequalities for proving the equivalence of all modes of convergence of $J_n$ for the $k$–nearest neighbor estimate under no conditions on the distribution of $(X, Y)$ other than the boundedness of $Y$. Thus, Stone's conditions on the relative sizes of $k$ and $n$ are strong enough to imply complete and almost sure convergence. Our other result is the strong universal convergence of $J_n$, that is, we can replace the boundedness assumption by the natural condition $\mathbf{E}|Y| < \infty$. Here we need the additional condition $k/\log n \to \infty$ on $k$.

Before we can state the main results, we have to take care of the messy problem of distance ties ($\|x - X_i\| = \|x - X_j\|$). The exponential inequality used here and in Devroye and Krzyżak (1989) is basically useful whenever the removal of one data point has a limited effect on the error. Also, our covering lemma requires some sort of duality that states that if $X_i$ is one of the near neighbors of $X_j$, then roughly speaking $X_j$ should be one of the near neighbors of $X_i$. Next we list three of the possible tie-breaking methods:

1. *Tie breaking by indices.* If $X_i$ and $X_j$ are equidistant from $x$, then $X_i$ is declared "closer" if $i < j$. This method has some undesirable properties. For example, if $X$ is monoatomic, then $X_1$ is the nearest neighbor of all $X_j$'s, $j > 1$, but $X_j$ is only the $(j - 1)$st nearest neighbor of $X_1$. The influence of $X_1$ in such a situation is too great, making the estimate very unstable and thus undesirable. In this case, Devroye and Györfi [(1985), Chapter 10] pointed out that, when $Y$ is not degenerate and $\varepsilon > 0$ is small enough,

$$\mathbf{P}\left\{ \int |\widetilde{m}_n(x) - m(x)| \mu(dx) > \varepsilon \right\} \geq \exp(-ck),$$

   for some $c > 0$. This is in contrast to Theorem 1(a) below, where $\widetilde{m}_n(x)$ is the $k$–nearest neighbor regression estimate defined by tie breaking by indices.

2. *Stone's tie breaking.* Stone (1977) introduced a nearest neighbor rule which is not a $k$–nearest neighbor rule in a strict sense, for his estimate, in general, uses more than $k$ neighbors. If we denote the distance of the $k$th nearest neighbor to $x$ by $R_n(x)$ (note that it is unique), then Stone's estimate is the following:

   $$\widehat{m}_n(x)$$

   (3)
   $$= \frac{1}{k}\left( \sum_{i:\, \|x - x_i\| < R_n(x)} Y_i + \frac{k - \#\{i:\, \|x - x_i\| < R_n(x)\}}{\#\{i:\, \|x - x_i\| = R_n(x)\}} \sum_{i:\, \|x - x_i\| = R_n(x)} Y_i \right).$$

3. *Tie breaking by randomization.* This is the method that we will consider. We assume that $(X, Z)$ is a random vector independent of the data, where $Z$ is independent of $X$ and uniformly distributed on $[0, 1]$. The latter assumption may be replaced by the weaker assumption that $Z$ has a density; however, as it is up to us to generate $Z$, we may as well pick a uniform random variable. We also artificially enlarge the data by introducing $Z_1, Z_2, \ldots, Z_n$, where the $Z_i$'s are i.i.d. uniform $[0,1]$ as well. Thus, each $(X_i, Z_i)$ is distributed as $(X, Z)$. The probability measure induced by $(X, Z)$ is denoted by $\nu$. Given $(x, z)$, we define

$$m_n(x, z) = \frac{1}{k}\sum_{i=1}^{k} Y_{(i)},$$

where

$$\left(X_{(1)}, Y_{(1)}\right), \ldots, \left(X_{(n)}, Y_{(n)}\right) = \left(X_{(1, x, z)}, Y_{(1, x, z)}\right), \ldots, \left(X_{(n, x, z)}, Y_{(n, x, z)}\right)$$

is a reordering of the data according to increasing values of $\|x - X_{(i)}\|$. In case of distance ties, we declare $(X_i, Z_i)$ closer to $(x, z)$ than $(X_j, Z_j)'$ provided that

$$|Z_i - z| \le |Z_j - z|.$$

The criterion is

$$J_n = \mathbf{E}\big\{|m_n(X, Z) - m(X)| \mid X_1, Z_1, Y_1, \ldots, X_n, Z_n, Y_n\big\}$$

$$= \int_0^1 \int |m_n(x, z) - m(x)| \mu(dx)\, dz = \int |m_n(x, z) - m(x)| \nu\big(d(x, z)\big).$$

The main difference between Stone's tie-breaking policy and the one based on randomization is that Stone's method takes into account all points whose distance to $x$ equals that of the $k$th nearest neighbor, while the method based on randomization picks one of these randomly and neglects the others. We will see in the proof of Theorem 1 that $\mathbf{E}J_n$ cannot be smaller than the expected $L_1$-error of Stone's estimate. It should be stressed that if $\mu$ has a density, then tie breaking is needed with zero probability and becomes therefore irrelevant.

**2. The equivalence theorem.** The purpose of this section is to prove the following result.

THEOREM 1. *Let $m_n(x, z)$ be the $k$–nearest neighbor estimate defined above. Then the following statements are equivalent:*

(a) *For every distribution of $(X, Y)$ with $|Y| \le M < \infty$ and $\varepsilon > 0$, there is a positive integer $n_0$ such that, for $n > n_0$,*

$$\mathbf{P}\{J_n > \varepsilon\} \le \exp\!\Big[-n\varepsilon^2/\big(8M^2\gamma_d^2\big)\Big],$$

*where the constant $\gamma_d$ is the minimal number of cones centered at the origin of angle $\pi/6$ that cover $\mathcal{R}^d$.*

(b) *For every distribution of $(X, Y)$ with $|Y| \le M < \infty$,*

$$J_n \to 0 \quad \text{with probability } 1 \text{ as } n \to \infty.$$

(c) *For every distribution of $(X, Y)$ with $|Y| \le M < \infty$,*

$$J_n \to 0 \quad \text{in probability as } n \to \infty.$$

(d) $\lim_{n \to \infty} k = \infty$ *and* $\lim_{n \to \infty} k/n = 0$.

REMARK 1 (A curiosity). It is interesting that we can find sequences $k$ for which $J_n \to 0$ almost surely for all distributions of $(X, Y)$ with bounded $Y$, yet $m_n$ does not tend to $m$ in the almost sure pointwise sense [take $k \sim \log\log\log(n)$, and note that $k/\log\log(n) \to \infty$ is necessary for the almost sure pointwise convergence of the kernel estimate whenever $X$ has a density and $m$ is twice continuously differentiable with $m'' \ne 0$; Devroye (1982)].

REMARK 2 (Necessity of the conditions). It is not true that, when $J_n \to 0$ in probability for one distribution of $(X, Y)$, the conditions (1) on $k$ follow: just consider the case that $Y = 0$ with probability 1. Of course, the implication is true for "most" distributions of $(X, Y)$.

REMARK 3 (General estimates). We will not consider smoothed versions of the $k$–nearest neighbor method here. For example, as in Devroye (1982), one might consider attaching weight $v_{ni}$ to the $i$th nearest neighbor, where $v_{n1} \geq v_{n2} \geq \cdots \geq v_{nn} \geq 0$ and the weights sum to 1 for every $n$. Such methods were first proposed by Royall (1966).

REMARK 4 (Other references). For other results on $k$–nearest neighbor convergence, see, for example, Collomb (1979, 1980), Mack (1981), Devroye (1978, 1981, 1982), Stute (1984) and Bhattacharya and Mack (1987).

REMARK 5 (Random $k$). If $k$ is replaced by a random variable $K$ that is independent of the data and satisfies $K/n \to 0$ and $K \to \infty$ almost surely, then $J_n \to 0$ almost surely. Such data-based choices can be obtained by splitting the data, for example.

REMARK 6 (Discrimination). The conditional probability of error, of the $k$–nearest neighbor rule in discrimination, given the data [Cover and Hart (1967)], converges completely and strongly to the Bayes probability of error as $n \to \infty$ for all distributions of the data whenever (1) holds. This result strengthens the universal weak convergence results of Stone (1977) and Devroye and Wagner (1980).

REMARK 7 ($L_p$-consistency). By the boundedness of $Y$ it is easy to see that the $L_p$-error

$$J_n^{(p)} = \left( \int_0^1 \int |m_n(x, z) - m(x)|^p \mu(dx) \, dz \right)^{1/p} \quad \text{for } 1 \leq p < \infty$$

converges to zero if and only if the $L_1$-error $J_n$ does; therefore the results of Theorem 1 remain valid for $L_p$-errors.

REMARK 8 (Inequalities). The inequality of Theorem 1(a) is less useful in practice, as it is only valid for $n > n_0$, where $n_0$ depends upon $\varepsilon$.

PROOF OF THEOREM 1. Clearly, (a) implies (b) and (b) implies (c). Part (c) implies that $\mathbf{E} J_n \to 0$; therefore, by Jensen's inequality,

$$\mathbf{E} J_n = \mathbf{E} \left( \int |m_n(x, z) - m(x)| \, \nu\big(d(x, z)\big) \right)$$

$$\geq \mathbf{E} \left( \int \left| \mathbf{E} \left( \int m_n(x, z) dz \, | X_1, Y_1, \ldots, X_n, Y_n \right) - m(x) \right| \mu(dx) \right)$$

$$= \mathbf{E} \left( \int |\widehat{m}_n(x) - m(x)| \mu(dx) \right) \to 0,$$

where $\widehat{m}_n$ is Stone's estimate defined by (3); but this implies (d) by the results of Stone (1977). The novelty in this paper is the proof that condition (d) implies

(a). We begin with an exponential inequality generalizing inequalities due to Hoeffding (1963). The generalization due to Azuma (1967) [see Stout (1974)] has led to interesting applications in combinatorics and the theory of random graphs [for a survey, see McDiarmid (1989)]. We have used it in density estimation [Devroye (1988, 1991)].

LEMMA 1 [McDiarmid (1989)].  *Let $X_1, \ldots, X_n$ be independent random variables taking values in a set $A$, and assume that $f \colon A^n \to R$ satisfies*

$$\sup_{\substack{x_1, \ldots, x_n \\ x'_1, \ldots, x'_n \in A}} |f(x_1, \ldots, x_n) - f(x_1, \ldots, x_{i-1}, x'_i, x_{i+1}, \ldots, x_n)| \le c_i, \qquad 1 \le i \le n.$$

*Then*

$$\mathbf{P}\{|f(X_1, \ldots, X_n) - \mathbf{E}f(X_1, \ldots, X_n)| \ge t\} \le 2 \exp\left(\frac{-2t^2}{\sum_{i=1}^n c_i^2}\right).$$

The other tool needed for our proof is exploiting some geometric properties of the "metric" defined by the tie-breaking rule. In order to make it more transparent, we recall Lemma 10.1 from Devroye and Györfi (1985), which was used in the proof of complete consistency of $k$–nearest neighbor estimates if $\mu$ has a density. Let $S_{x,r}$ and $\overline{S}_{x,r}$ denote the open and closed balls of radius $r$ centered at $x$, respectively.

LEMMA 2 [Devroye and Györfi (1985)].  *Let $\mu$ be an absolutely continuous probability measure on $\mathcal{R}^d$. Define*

$$B_a(x) = \left\{x' \colon \mu\big(\overline{S}_{x', \|x - x'\|}\big) \le a\right\}.$$

*Then, for all $x \in \mathcal{R}^d$,*

$$\mu\big(B_a(x)\big) \le \gamma_d a.$$

Since Devroye and Györfi assumed the existence of a density, they did not have to worry about tie breaking. In order to generalize Lemma 2 to our case, we need some notation. For $x \in \mathcal{R}^d$ let $C(x) \subset \mathcal{R}^d$ be a cone of angle $\pi/6$ centered at $x$. The cone consists of all $y$ with the property that either $y = x$ or angle$(y - x, s) \le \pi/6$, where $s$ is a fixed direction. If $y, y' \in C(x)$ and $\|x - y\| < \|x - y'\|$, then $\|y - y'\| < \|x - y'\|$. Furthermore, if $\|x - y\| \le \|x - y'\|$, then $\|y - y'\| \le \|x - y'\|$. This follows from a simple geometric argument in the vector space spanned by $x$, $y$ and $y'$.

For $(x, z) \in \mathcal{R}^d \times [0, 1]$ define $C_0(x, z), C_1(x, z), S_{(x,z),(r,b)} \subset \mathcal{R}^d \times [0, 1]$ as

$$C_0(x, z) = C(x) \times [0, z],$$
$$C_1(x, z) = C(x) \times [z, 1]$$

and

$$S_{(x,z),(r,b)} = S_{x,r} \times [0,1] \cup \left\{(\widehat{x},\widehat{z}): \|\widehat{x} - x\| = r, |\widehat{z} - z| \leq b\right\}.$$

Clearly, $\mathcal{R}^d \times [0,1]$ can be covered by $2\gamma_d$ sets of type $C_0(x,z)$ and $C_1(x,z)$. The property that we need is the following.

LEMMA 3. *Let*
$$B_a(x,z) = \left\{(x',z'): \nu\left(S_{(x',z'),(\|x-x'\|,|z-z'|)}\right) \leq a\right\}.$$
*Then for all $(x,z) \in \mathcal{R}^d \times [0,1]$*

$$\nu\left(B_a(x,z)\right) \leq 2\gamma_d a.$$

First of all, we prove a covering lemma, the key property of the "cones" $C_0(x,z)$ and $C_1(x,z)$.

LEMMA 4. *If $(x',z') \in C_0(x,z)$, then*

$$C_0(x,z) \cap S_{(x,z),(\|x-x'\|,|z-z'|)} \subset S_{(x',z'),(\|x-x'\|,|z-z'|)},$$

*and if $(x',z') \in C_1(x,z)$, then*

$$C_1(x,z) \cap S_{(x,z),(\|x-x'\|,|z-z'|)} \subset S_{(x',z'),(\|x-x'\|,|z-z'|)}.$$

PROOF. Because of symmetry it is enough to prove one of the statements. We have to show that $(\widehat{x},\widehat{z}) \in C_0(x,z) \cap S_{(x,z),(\|x-x'\|,|z-z'|)}$ implies $(\widehat{x},\widehat{z}) \in S_{(x',z'),(\|x-x'\|,|z-z'|)}$.

If $\widehat{x} \in C(x) \cap S_{x,\|x-x'\|}$, then from the well-known property of the cone $\widehat{x} \in S_{x',\|x-x'\|}$ follows, so it is enough to deal with pairs $(\widehat{x},\widehat{z})$ where $\|\widehat{x} - x\| = \|x - x'\|$. Since $\widehat{x} \in C(x)$, the only case when $\widehat{x} \notin S_{x',\|x-x'\|}$ is if

$$\|\widehat{x} - x\| = \|\widehat{x} - x'\| = \|x' - x\|.$$

Denote the set of such $\widehat{x}$'s by $H$. Thus, it is enough to deal with pairs in $H \times [0,1]$. Intersecting this set with the left- and right-hand sides of the statement, we get

$$H \times [0,1] \cap C_0(x,z) \cap S_{(x,z),(\|x-x'\|,|z-z'|)}$$
$$= H \times \left([0,z] \cap \{\widehat{z}: |\widehat{z} - z| \leq |z - z'|\}\right) = H \times [z',z]$$

and

$$H \times [0,1] \cap S_{(x',z'),(\|x-x'\|,|z-z'|)} = H \times \left(\{\widehat{z}: |\widehat{z} - z'| \leq |z - z'|\}\right),$$

respectively. Clearly, however,

$$[z',z] \subset \{\widehat{z}: |\widehat{z} - z'| \leq |z - z'|\},$$

which completes the proof. □

PROOF OF LEMMA 3.   The proof is similar to that of Lemma 2. Let $C_1, \ldots,$ $C_{2\gamma_d}$ be a collection of sets of form $C_0(x,z)$ and $C_1(x,z)$ that cover $\mathcal{R}^d \times [0,1]$. Then

$$\nu\big(B_a(x,z)\big) \leq \sum_{s=1}^{2\gamma_d} \nu\big(C_s \cap B_a(x,z)\big).$$

Let $(x',z') \in C_s \cap B_a(x,z)$. Then from Lemma 4 we have

$$\nu\big(C_s \cap S_{(x,z),(\|x-x'\|,\,|z-z'|)} \cap B_a(x,z)\big) \leq \nu(S_{(x',z'),(\|x-x'\|,\,|z-z'|)}) \leq a,$$

where we used the fact that $(x',z') \in B_a(x,z)$. Since $(x',z')$ was arbitrary,

$$\nu\big(C_s \cap B_a(x,z)\big) \leq a,$$

which completes the proof of the lemma.   □

Now we are equipped to prove that (d) implies (a) in Theorem 1. Set $r_n = r_n(x,z)$ and $b_n = b_n(x,z)$ to satisfy

$$\nu\big(S_{(x,z),(r_n,b_n)}\big) = \frac{k}{n}.$$

Note that the solution always exists, by the absolute continuity of the distribution of $Z$ and its independence from $X$. Also define

$$m_n^*(x,z) = \frac{1}{k} \sum_{j=1}^{n} Y_j I_{\{(X_j,Z_j) \in S_{(x,z),(r_n,b_n)}\}}.$$

Obviously,

$$(4) \qquad \begin{aligned} |m(x) - m_n(x,z)| &\leq \big|m(x) - \mathbf{E}\big(m_n^*(x,z)\big)\big| \\ &\quad + \big|\mathbf{E}\big(m_n^*(x,z)\big) - m_n^*(x,z)\big| + |m_n^*(x,z) - m_n(x,z)|. \end{aligned}$$

The first term on the right-hand side is a deterministic "bias"-type term, whose integral will be shown to converge to zero. The second and third terms are random; they can be considered as "variation" terms. We will obtain exponential probability inequalities for these terms that are valid for large $n$'s.

The condition $k/n \to 0$ implies that $r_n(x,z) \to 0$, so for the first term we have, by Lebesgue's density theorem [see Wheeden and Zygmund (1977)], that

$$\begin{aligned} \mathbf{E}\big(m_n^*(x,z)\big) &= \frac{1}{\nu\big(S_{(x,z),(r_n,b_n)}\big)} \int_{S_{(x,z),(r_n,b_n)}} \mathbf{E}\big(Y \mid (X,Z) = (x',z')\big)\nu\big(d(x',z')\big) \\ &\to \mathbf{E}\big(Y \mid (X,Z) = (x,z)\big) = m(x) \end{aligned}$$

for almost all $x$ mod $\mu$. By the boundedness of $Y$, the dominated convergence theorem implies that

$$\int \big|m(x) - \mathbf{E}\big(m_n^*(x,z)\big)\big| \nu\big(d(x,z)\big) \to 0.$$

Turning to the second term in (4), first we get an exponential bound for

$$\left| \int \left| \mathbf{E}\big(m_n^*(x,z)\big) - m_n^*(x,z) \right| \nu\big(d(x,z)\big) - \mathbf{E} \int \left| \mathbf{E}\big(m_n^*(x,z)\big) - m_n^*(x,z) \right| \nu\big(d(x,z)\big) \right|$$

by Lemma 1. Fix the data and replace $(x_i, z_i, y_i)$ by $(\widehat{x}_i, \widehat{z}_i, \widehat{y}_i)$, changing the value of $m_n^*(x,z)$ to $m_{ni}^*(x,z)$. Then

$$\left| \int \left| \mathbf{E}\big(m_n^*(x,z)\big) - m_n^*(x,z) \right| \nu\big(d(x,z)\big) - \int \left| \mathbf{E}\big(m_n^*(x,z)\big) - m_{ni}^*(x,z) \right| \nu\big(d(x,z)\big) \right|$$

$$\leq \int \left| m_n^*(x,z) - m_{ni}^*(x,z) \right| \nu\big(d(x,z)\big);$$

but $|m_n^*(x,z) - m_{ni}^*(x,z)|$ is bounded by $2M/k$ and can differ from zero only if $(x_i, z_i) \in S_{(x,z),(r_n, b_n)}$ or $(\widehat{x}_i, \widehat{z}_i) \in S_{(x,z),(r_n, b_n)}$. Observe that $(x_i, z_i) \in S_{(x,z),(r_n, b_n)}$ if and only if $\nu(S_{(x,z),(\|x - x_i\|, |z - z_i|)}) \leq k/n$. However, the measure of such $(x,z)$ pairs is bounded by $\gamma_d k/n$, by Lemma 3; therefore,

$$\sup_{x_1,y_1,z_1,\ldots,x_n,y_n,z_n,\widehat{x}_i,\widehat{z}_i,\widehat{y}_i} \int \left| m_n^*(x,z) - m_{ni}^*(x,z) \right| \nu\big(d(x,z)\big) \leq \frac{2M}{k}\frac{\gamma_d k}{n} = \frac{2M\gamma_d}{n}$$

and, by Lemma 1,

$$(5) \qquad \begin{aligned} \mathbf{P}\bigg\{ &\left| \int |\mathbf{E}m_n^*(x,z) - m_n^*(x,z)| \nu\big(d(x,z)\big) \right. \\ &\qquad \left. - \mathbf{E} \int |\mathbf{E}m_n^*(x,z) - m_n^*(x,z)| \nu\big(d(x,z)\big) \right| > \varepsilon \bigg\} \\ &\leq 2e^{-n\varepsilon^2/(2M^2\gamma_d^2)}. \end{aligned}$$

So we have to show that

$$\mathbf{E} \int |\mathbf{E}m_n^*(x,z) - m_n^*(x,z)| \nu\big(d(x,z)\big) \to 0.$$

However, using the Cauchy–Schwarz inequality, we have

$$\mathbf{E} \int |\mathbf{E}m_n^*(x,z) - m_n^*(x,z)| \nu\big(d(x,z)\big)$$

$$\leq \int \sqrt{\mathbf{E}|\mathbf{E}m_n^*(x,z) - m_n^*(x,z)|^2} \, \nu\big(d(x,z)\big)$$

$$= \int \sqrt{\frac{1}{k^2} n \operatorname{Var}\big(YI_{\{(X,Z) \in S_{(x,z),(r_n,b_n)}\}}\big)} \, \nu\big(d(x,z)\big)$$

$$\leq \int \sqrt{\frac{M^2}{k^2} n \nu\big(S_{(x,z),(r_n,b_n)}\big)} \, \nu\big(d(x,z)\big)$$

$$= \int \sqrt{\frac{nM^2}{k^2}\frac{k}{n}} \, \mu(dx)$$

$$= \sqrt{\frac{M^2}{k}} \to 0.$$

Finally, denoting $R_n = \|X_{(k)} - x\|$ and $B_n = |Z_{(k)} - z|$, write the third term in (4) as

$$
\begin{aligned}
&|m_n^*(x, z) - m_n(x, z)| \\
&= \frac{1}{k} \left| \sum_{j=1}^{n} Y_j I_{\{(X_j, Z_j) \in S_{(x, z), (r_n, b_n)}\}} - \sum_{j=1}^{n} Y_j I_{\{(X_j, Z_j) \in S_{(x, z), (R_n, B_n)}\}} \right| \\
&\leq \frac{M}{k} \sum_{j=1}^{n} \left| I_{\{(X_j, Z_j) \in S_{(x, z), (r_n, b_n)}\}} - I_{\{(X_j, Z_j) \in S_{(x, z), (R_n, B_n)}\}} \right| \\
&= M \left| \frac{1}{k} \sum_{j=1}^{n} I_{\{(X_j, Z_j) \in S_{(x, z), (r_n, b_n)}\}} - 1 \right| \\
&= M |\widehat{m}_n^*(x, z) - \mathbf{E}\widehat{m}_n^*(x, z)|,
\end{aligned}
$$

where $\widehat{m}_n^*$ is defined as $m_n^*$ with $Y$ replaced by the constant random variable $\widehat{Y} = 1$. Therefore the bound of (5) applies for the third term, too, and the proof is complete. □

**3. Strong universal consistency.** In this section we demonstrate that the $k$–nearest neighbor regression estimate is consistent even if $Y$ is not bounded, if $k$ is chosen to satisfy $k/\log(n) \to \infty$ and $k/n \to 0$. More precisely, we prove the following theorem.

THEOREM 2.    *If*

$$
\lim_{n \to \infty} k/\log(n) = \infty \quad and \quad \lim_{n \to \infty} k/n = 0,
$$

*then $J_n \to 0$ with probability 1 for all distributions of $(X, Y)$ satisfying $\mathbf{E}|Y| < \infty$.*

Györfi (1991) gave conditions for the strong universal consistency of a regression estimate. Translating his result to our case, we get the following lemma.

LEMMA 5 [Györfi (1991), Theorem 2].    *Consider the $k$–nearest neighbor regression estimate $m_n(x, z)$. Then the $L_1$-error of the estimate $J_n$ converges to zero almost surely for all distributions of $(X, Y)$ satisfying $\mathbf{E}|Y| < \infty$ if the following two conditions are satisfied:*

(a) *$J_n \to 0$ almost surely for all distributions of $(X, Y)$ with bounded $Y$.*

(b) *There is a constant $c > 0$ such that, for all distributions of $(X, Y)$ satisfying $\mathbf{E}|Y| < \infty$,*

$$
\limsup_{n \to \infty} \frac{1}{k} \sum_{i=1}^{k} \int_0^1 \int |Y_{(i, x, z)}| \mu(dx) \, dz \leq c\mathbf{E}|Y| \quad a.s.
$$

Clearly, condition (a) is satisfied by Theorem 1, so we only have to check (b). In order to do so, we need some notation. Let $A_i$ be the collection of all $(x, z)$

that are such that $(X_i, Z_i)$ is one of its $k$ nearest neighbors. Here, we use some geometric arguments similar to those in the proof of Theorem 1. Similarly, let us define a cone $C(x, \theta, s)$, where $x$ defines the top of the cone, $s$ is a vector indicating a direction in $\mathcal{R}^d$ and $\theta \in (0, \pi)$ is an angle. The cone consists of all $y$ with the property that either $y = x$ or angle$(y - x, s) \leq \theta$. For any fixed $\theta$, there exists a finite collection $S$ of directions such that

$$\bigcup_{s \in S} C(x, \theta, s) = \mathcal{R}^d$$

regardless of how $x$ is picked. The cardinality of this set is denoted by $|S|$ and depends upon both $\theta$ and $d$. If $\theta < \pi/6$ and if $y, y' \in C(x, \theta, s)$ and $\|x - y\| < \|x - y'\|$, then $\|y - y'\| < \|x - y'\|$. Furthermore, if $\|x - y\| \leq \|x - y'\|$, then $\|y - y'\| \leq \|x - y'\|$. We fix $\theta \in (0, \pi/6)$ and $S$ as indicated above. In the space $\mathcal{R}^d \times [0, 1]$, define the sets

$$C_{i,s} = C(X_i, \theta, s) \times [0, 1].$$

Let $B_{i,s}$ be the subset of $C_{i,s}$ consisting of all $(x, z)$ that are among the $k$ nearest neighbors of $(X_i, Z_i)$ in the set

$$\{(X_1, Z_1), \dots, (X_{i-1}, Z_{i-1}), (X_{i+1}, Z_{i+1}), \dots, (X_n, Z_n), (x, z)\} \cap C_{i,s}$$

when distance tie breaking is done in the described fashion. [If $C_{i,s}$ contains fewer than $k - 1$ of the $(X_j, Z_j)$ pairs $i \neq j$, then $B_{i,s} = C_{i,s}$.] Equivalently, $B_{i,s}$ is the subset of $C_{i,s}$ consisting of all $(x, z)$ that are closer to $(X_i, Z_i)$ than the $k$th nearest neighbor of $(X_i, Z_i)$ in $C_{i,s}$, when distance tie breaking is done by the described fashion.

LEMMA 6. *If $(x, z) \in A_i$, then $(x, z) \in \cup_{s \in S} B_{i,s}$, and thus*

$$\nu(A_i) \leq \sum_{s \in S} \nu(B_{i,s}).$$

PROOF. To prove this claim, take $(x, z) \in A_i$. Then locate an $s \in S$ for which $(x, z) \in C_{i,s}$. We have to show that $(x, z) \in B_{i,s}$ to conclude the proof. Thus, we need to show that $(x, z)$ is one of the $k$ nearest neighbors of $(X_i, Z_i)$ in the set

$$\{(X_1, Z_1), \dots, (X_{i-1}, Z_{i-1}), (X_{i+1}, Z_{i+1}), \dots, (X_n, Z_n), (x, z)\} \cap C_{i,s}$$

when distance tie breaking is done appropriately. Take $(X_j, Z_j)$ closer to $(X_i, Z_i)$ than $(x, z)$ in $C_{i,s}$. If $\|X_j - X_i\| < \|x - X_i\|$, we recall that by the property of our cones $\|x - X_j\| < \|x - X_i\|$, and thus $(X_j, Z_j)$ is one of the $k - 1$ nearest neighbors of $(x, z)$ in $\mathcal{R}^d$. If on the other hand $\|X_j - X_i\| = \|x - X_i\|$, and $z$ is further from $Z_i$ than $Z_j$, then by the property of the cone, $\|x - X_j\| < \|x - X_i\|$, which shows again that $(X_j, Z_j)$ is one of the $k - 1$ nearest neighbors of $(x, z)$ in $\mathcal{R}^d$. This shows that in $C_{i,s}$ there are at most $k - 1$ points $(X_j, Z_j)$ closer to $(X_i, Z_i)$ than $(x, z)$.

Thus, with the same tie-breaking policy, $(x, z)$ is one of the $k$ nearest neighbors of $(X_i, Z_i)$ in the set

$$\big\{(X_1, Z_1), \ldots, (X_{i-1}, Z_{i-1}), (X_{i+1}, Z_{i+1}), \ldots, (X_n, Z_n), (x, z)\big\} \cap C_{i,s}.$$

This concludes the proof of the claim. □

LEMMA 7 (An inequality for binomial random variables). *Let $B$ be a binomial random variable with parameters $n$ and $p$. Then*

$$\mathbf{P}\{B > \varepsilon\} \leq \exp\Big[\varepsilon - np - \varepsilon \log(\varepsilon/np)\Big], \qquad \varepsilon > np,$$

$$\mathbf{P}\{B < \varepsilon\} \leq \exp\Big[\varepsilon - np - \varepsilon \log(\varepsilon/n)\Big], \qquad \varepsilon < np.$$

PROOF. We proceed by Chernoff's exponential bounding method [Chernoff (1952)]. In particular, for arbitrary $\lambda > 0$,

$$\begin{aligned}
\mathbf{P}\{B > \varepsilon\} &\leq \mathbf{E}\{\exp(\lambda B - \lambda \varepsilon)\} \\
&= \exp(-\lambda \varepsilon)\big((\exp \lambda)p + 1 - p\big)^n \\
&\leq \exp\Big[-\lambda \varepsilon + np\big((\exp \lambda) - 1\big)\Big].
\end{aligned}$$

The right-hand side is minimal for $\lambda = \log(\varepsilon/np)$. Resubstitution of this value gives the first bound. The proof of the other bound is similar. □

LEMMA 8 [Property of $\nu(B_{i,s})$]. *If $k/\log(n) \to \infty$ and $k/n \to 0$, then*

$$\limsup_{n \to \infty} \frac{n}{k} \max_i \nu(B_{i,s}) \leq 2 \quad a.s.$$

PROOF. We prove that, for every $s \in S$,

$$\sum_n \mathbf{P}\bigg\{\frac{n}{k} \max_i \nu(B_{i,s}) > 2\bigg\} < \infty.$$

In order to do this we give a bound for

$$\mathbf{P}\{\nu(B_{i,s}) > \varepsilon \mid X_i, Z_i\}.$$

If $\nu(C_{i,s}) \leq \varepsilon$, then, since $B_{i,s} \subseteq C_{i,s}$, we have $\mathbf{P}\{\nu(B_{i,s}) > \varepsilon \mid X_i, Z_i\} = 0$; therefore we assume that $\nu(C_{i,s}) > \varepsilon$. Fix $X_i$ and $Z_i$. The distance-ordering and tie-breaking method induces a total ordering of all $(x, z)$ with respect to closeness to $(X_i, Z_i)$. Find a pair $(x, z) \in C_{i,s}$ such that if $B_\varepsilon$ is the collection of all $(x', z') \in C_{i,s}$ that are nearer to $(X_i, Z_i)$ than $(x, z)$, then $\nu(B_\varepsilon) = \varepsilon$. By our method of tie breaking, such a pair $(x, z)$ exists. We have the following dual relationship:

$$\begin{aligned}
&\mathbf{P}\{\nu(B_{i,s}) > \varepsilon \mid X_i, Z_i\} \\
&\quad = \mathbf{P}\{B_\varepsilon \text{ captures fewer than } k \text{ of the points } (X_j, Z_j), j \neq i \mid X_i, Z_i\}.
\end{aligned}$$

However, if $B$ is a binomial $(n - 1, \varepsilon)$ random variable, then the last probability is equal to

$$\mathbf{P}\{B < k\} \leq \exp\left[k - (n - 1)\varepsilon - k\log\left(\frac{k}{(n-1)\varepsilon}\right)\right] \quad \text{if } k < (n - 1)\varepsilon.$$

Finally, with $\varepsilon = 2k/n$ we have $k < (n - 1)\varepsilon$; therefore

$$\mathbf{P}\left\{\max_{1 \leq i \leq n} \nu(B_{i,s}) > \varepsilon\right\} \leq n\mathbf{P}\{\nu(B_{1s}) > \varepsilon\}$$

$$\leq n\exp\left[k - (n - 1)\varepsilon - k\log\left(\frac{k}{(n-1)\varepsilon}\right)\right]$$

$$= n\exp\left[k - \frac{2k(n - 1)}{n} - k\log\left(\frac{n}{2(n-1)}\right)\right]$$

$$\leq n\exp\left(-k + \frac{2k}{n} + k\log 2\right),$$

which is summable in $n$ when $k \geq [2/(1 - \log 2)]\log n$. $\square$

PROOF OF THEOREM 2. By Lemma 5 and Theorem 1, it is enough to prove that there is a constant $c > 0$ such that

$$\limsup_{n \to \infty} \frac{1}{k}\sum_{i=1}^{k}\int_{0}^{1}\int |Y_{(i,x,z)}|\mu(dx)\,dz \leq c\mathbf{E}|Y| \quad \text{a.s.}$$

Observe that

$$\frac{1}{k}\sum_{i=1}^{k}\int_{0}^{1}\int |Y_{(i,x,z)}|\mu(dx)\,dz = \frac{1}{k}\sum_{i=1}^{n}|Y_{(i)}|\nu(A_i) \leq \frac{1}{n}\sum_{i=1}^{n}|Y_{(i)}|\left(\frac{n}{k}\max_i \nu(A_i)\right).$$

If we can show that

$$(6) \qquad \limsup_{n \to \infty} \frac{n}{k}\max_i \nu(A_i) \leq c \quad \text{a.s.},$$

for some constant $c > 0$, then, by the law of large numbers,

$$\limsup_{n \to \infty} \frac{1}{n}\sum_{i=1}^{n}|Y_{(i)}|\left(\frac{n}{k}\max_i \nu(A_i)\right) \leq \limsup_{n \to \infty} c\frac{1}{n}\sum_{i=1}^{n}|Y_{(i)}| = c\mathbf{E}|Y| \quad \text{a.s.},$$

so we have to prove (6). However, by Lemma 6,

$$\nu(A_i) \leq \sum_{s \in S}\nu(B_{i,s}),$$

therefore, Lemma 8 implies that (6) is satisfied with $c = 2|S|$, so the proof of the theorem is complete. $\square$

## REFERENCES

AZUMA, K. (1967). Weighted sums of certain dependent random variables. *Tôhoku Math. J.* **37** 357–367.

BECK, J. (1979). The exponential rate of convergence of error for $k_n$ NN nonparametric regression and decision. *Problems Control Inform. Theory* **8** 303–311.

BHATTACHARYA, P. K. and MACK, Y. P. (1987). Weak convergence of $k$-NN density and regression estimators with varying $k$ and applications. *Ann. Statist.* **15** 976–994.

CHERNOFF, H. (1952). A measure of asymptotic efficiency of tests of a hypothesis based on the sum of observations. *Ann. Math. Statist.* **23** 493–507.

COLLOMB, G. (1979). Estimation de la regression par la méthode des $k$ points les plus proches: propriétés de convergence ponctuelle. *C. R. Acad. Sci. Paris* **289** 245–247.

COLLOMB, G. (1980). Estimation de la regression par la méthode des $k$ points les plus proches avec noyau. *Statistique non Paramétrique Asymptotique. Lecture Notes in Math.* **821** 159–175. Springer, Berlin.

COLLOMB, G. (1981). Estimation non paramétrique de la regression: revue bibliographique. *Internat. Statist. Rev.* **49** 75–93.

COLLOMB, G. (1985). Nonparametric regression: an up-to-date bibliography. *Statistics* **16** 300–324.

COVER, T. M. (1968). Estimation by the nearest neighbor rule. *IEEE Trans. Inform. Theory* **IT-14** 50–55.

COVER, T. M. and HART, P. E. (1967). Nearest neighbor pattern classification. *IEEE Trans. Inform. Theory* **IT-13** 21–27.

DEVROYE, L. (1978). A universal $k$–nearest neighbor procedure in discrimination. In *Proceedings of the 1978 IEEE Computer Society Conference on Pattern Recognition and Image Processing* 142–147. IEEE, New York.

DEVROYE, L. (1981). On the almost everywhere convergence of nonparametric regression function estimates. *Ann. Statist.* **9** 1310–1319.

DEVROYE, L. (1982). Necessary and sufficient conditions for the almost everywhere convergence of nearest neighbor regression function estimates. *Z. Wahrsch. Verw. Gebiete* **61** 467–481.

DEVROYE, L. (1983). The equivalence of weak, strong and complete convergence in L1 for kernel density estimates. *Ann. Statist.* **11** 896–904.

DEVROYE, L. (1988). The kernel estimate is relatively stable. *Probab. Theory Related Fields* **77** 521–536.

DEVROYE, L. (1991). Exponential inequalities in nonparametric estimation. In *Nonparametric Functional Estimation* (G. Roussas, ed.) 31–44. Springer, Berlin.

DEVROYE, L. and GYÖRFI, L. (1983). Distribution-free exponential bound on the $L_1$ error of partitioning estimates of a regression function. In *Proceedings of the Fourth Pannonian Symposium on Mathematical Statistics* (F. Konecny, J. Mogyorodi and W. Wertz, eds.) 67–76. Akademiai Kiado, Budapest.

DEVROYE, L. and GYÖRFI, L. (1985). *Nonparametric Density Estimation: The $L_1$ View.* Wiley, New York.

DEVROYE, L. and KRZYŻAK, A. (1989). An equivalence theorem for L1 convergence of the kernel regression estimate. *J. Statist. Plann. Inference* **23** 71–82.

DEVROYE, L. and WAGNER, T. J. (1980). Distribution-free consistency results in nonparametric discrimination and regression function estimation. *Ann. Statist.* **8** 231–239.

GLICK, N. (1974). Consistency conditions for probability estimators and integrals of density estimators. *Utilitas Math.* **6** 61–74.

GREBLICKI, W., KRZYŻAK, A. and PAWLAK, M. (1984). Distribution-free pointwise consistency of kernel regression estimate. *Ann. Statist.* **12** 1570–1575.

GYÖRFI, L. (1981). Recent results on nonparametric regression estimate and multiple classification. *Problems Control Inform. Theory* **10** 43–52.

GYÖRFI, L. (1991). Universal consistencies of a regression estimate for unbounded regression functions. In *Nonparametric Functional Estimation* (G. Roussas, ed.) 329–338. Springer, Berlin.

HART, P. E. (1968). The condensed nearest neighbor rule. *IEEE Trans. Inform. Theory* **IT-14** 515–516.

HOEFFDING, W. (1963). Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.* **58** 13–30.

KRZYŻAK, A. (1986). The rates of convergence of kernel regression estimates and classification rules. *IEEE Trans. Inform. Theory* **IT-32** 668–679.

KRZYŻAK, A. and PAWLAK, M. (1984). Distribution-free consistency of a nonparametric kernel regression estimate and classification. *IEEE Trans. Inform. Theory* **IT-30** 78–81.

MACK, Y. P. (1981). Local properties of $k$–nearest neighbor regression estimates. *SIAM Journal on Algebraic and Discrete Methods* **2** 311–323.

McDIARMID, C. (1989). On the method of bounded differences. In *Surveys in Combinatorics 1989. London Mathematical Society Lecture Notes Series* **141** 148–188, Cambridge Univ. Press.

NADARAYA, E. A. (1964). On estimating regression. *Theory Probab. Appl.* **9** 141–142.

NADARAYA, E. A. (1970). Remarks on nonparametric estimates for density functions and regression curves. *Theory Probab. Appl.* **15** 134–137.

ROYALL, R. M. (1966). A class of nonparametric estimators of a smooth regression function. Ph.D. dissertation, Stanford Univ.

SPIEGELMAN, C. and SACKS, J. (1980). Consistent window estimation in nonparametric regression. *Ann. Statist.* **8** 240–246.

STONE, C. J. (1977). Consistent nonparametric regression (with discussion). *Ann. Statist.* **5** 595–645.

STOUT, W. F. (1974). *Almost Sure Convergence*. Academic, New York.

STUTE, W. (1984). Asymptotic normality of nearest neighbor regression function estimates. *Ann. Statist.* **12** 917–926.

WATSON, G. S. (1964). Smooth regression analysis. *Sankhyā Ser. A* **26** 359–372.

WHEEDEN, R. L. and ZYGMUND, A. (1977). *Measure and Integral*. Dekker, New York.

ZHAO, L. C. (1987). Exponential bounds of mean error for the nearest neighbor estimates of regression functions. *J. Multivariate Anal.* **21** 168–178.

LUC DEVROYE
SCHOOL OF COMPUTER SCIENCE
McGILL UNIVERSITY
MONTREAL
CANADA H3A 2A7

ADAM KRZYŻAK
DEPARTMENT OF COMPUTER SCIENCE
CONCORDIA UNIVERSITY
1455 DE MAISONNEUVE WEST
MONTREAL
CANADA H3G 1M8

LÁSZLÓ GYÖRFI
GÁBOR LUGOSI
DEPARTMENT OF MATHEMATICS
TECHNICAL UNIVERSITY OF BUDAPEST
1521 STOCZEK U. 2
BUDAPEST
HUNGARY