

A COMPARISON OF CROSS-VALIDATION TECHNIQUES IN DENSITY ESTIMATION¹

By J. S. MARRON

University of North Carolina, Chapel Hill

In the setting of nonparametric multivariate density estimation, theorems are established which allow a comparison of the Kullback–Leibler and the least-squares cross-validation methods of smoothing parameter selection. The family of delta sequence estimators (including kernel, orthogonal series, histogram and histospline estimators) is considered. These theorems also show that either type of cross validation can be used to compare different estimators (e.g., kernel versus orthogonal series).

1. Introduction. Consider the problem of trying to estimate a d -dimensional probability density function, $f(x)$, using a random sample, X_1, \dots, X_n , from f . Most proposed estimators of f depend on a “smoothing parameter,” say $\lambda \in \mathbb{R}^+$, whose selection is crucial to the performance of the estimator.

In this paper, for the large class of delta sequence estimators, theorems are obtained which allow comparison of two smoothing parameter selectors which are known to be asymptotically optimal. An important consequence of these results is that either smoothing parameter selector may be used for a data based comparison of two density estimators, for example, kernel versus orthogonal series. Another attractive feature of these results is that they are set in a quite general framework, special cases of which provide simpler proofs of several recent asymptotic optimality results.

In Sections 2 and 3 the family of delta sequence estimators and the smoothing parameter selectors are given. The theorems are stated in Section 4, with some remarks in Section 5. The rest of the paper consists of proofs.

2. Delta sequence estimators. A delta sequence density estimator, as studied by Watson and Leadbetter (1965), Földes and Révész (1974) and Walter and Blum (1979), is an estimator which can be written in the form

$$\hat{f}_\lambda(x) = n^{-1} \sum_{i=1}^n \delta_\lambda(x, X_i),$$

where the function $\delta_\lambda(x, y)$ is indexed by the smoothing parameter $\lambda \in \mathbb{R}^+$.

Received January 1985; revised March 1986.

¹Research partially supported by ONR contract N00014-81-K-0373 and by Sonderforschungsbereich 123 of the Deutsche Forschungsgemeinschaft.

AMS 1980 *subject classifications*. Primary 62G05; secondary 62G20.

Key words and phrases. Cross validation, smoothing parameter selection, choice of nonparametric estimators, bandwidth selection.

Examples include:

Kernel estimators. Given a kernel function $K(x)$, let

$$\delta_\lambda(x, X_i) = \lambda K(\lambda^{1/d}(x - X_i)).$$

Histogram estimators. Let B denote a bounded subset of \mathbb{R}^d and suppose that A_1, \dots, A_λ form a partition of B . For $k = 1, \dots, \lambda$, let $1_k(x)$ denote the indicator of A_k and let μ denote Lebesgue measure. Define, for $\lambda \in \mathbb{Z}^+$,

$$\delta_\lambda(x, X_i) = \sum_{k=1}^{\lambda} \mu(A_k)^{-1} 1_k(x) 1_k(X_i).$$

Orthogonal series estimators. Given a nonnegative weight function w and a sequence of functions $\{\psi_k(x)\}$ which is orthonormal and complete with respect to the inner product

$$\int \psi_k(x) \psi_{k'}(x) w(x) dx,$$

define, for $\lambda \in \mathbb{Z}^+$,

$$\delta_\lambda(x, X_i) = \sum_{k=1}^{\lambda} \psi_k(x) \psi_k(X_i) w(x).$$

See Walter and Blum (1979) for a rather extensive list of other delta sequence density estimators.

3. Smoothing parameter selectors. The two methods of choosing the smoothing parameter λ that are discussed in this paper are the Kullback–Leibler and the least-squares methods of cross validation. Both make use of the leave-one-out estimators:

$$\hat{f}_{\lambda, j^-}(x) = (n - 1)^{-1} \sum_{i \neq j} \delta_\lambda(x, X_i), \quad j = 1, \dots, n.$$

The Kullback–Leibler (also known as pseudo-likelihood) method first appeared in Habbema, Hermans and van den Broeck (1974) and was modified in Marron (1985). This involves choosing λ to maximize

$$\text{KL}(\lambda) = \prod_{j=1}^n \left[\hat{f}_{\lambda, j^-}^+(X_j)^{u(X_j)} e^{-\hat{p}(\lambda)} \right],$$

where $\hat{f}_{\lambda, j^-}^+(x)$ is the positive part of $\hat{f}_{\lambda, j^-}(x)$,

$$\hat{f}_{\lambda, j^-}^+(x) = \hat{f}_{\lambda, j^-}(x) \vee 0,$$

where $u(x)$ is a nonnegative weight function which is supported on a set where f is bounded above 0 (for example the indicator of such a set), and where

$$\hat{p}(\lambda) = \int \hat{f}_\lambda(x) u(x) dx.$$

The least-squares method was introduced by Rudemo (1982) and Bowman (1984). This involves choosing λ to minimize

$$\text{LS}(\lambda) = \int \hat{f}_\lambda(x)^2 w(x) dx - 2n^{-1} \sum_{j=1}^n \hat{f}_{\lambda, j^-}(x) w(X_j),$$

where $w(x)$ is a nonnegative weight function.

For purposes of comparison of these two smoothing parameter selectors, the natural connection between the weight functions u and w will be seen to be

$$u(x) = w(x)f(x).$$

4. Theorems. In this section it will be demonstrated that choosing λ by the methods in the last section is, in a strong sense, asymptotically equivalent to minimizing the following distances:

Average square error.

$$d_A(\hat{f}_\lambda, f) = n^{-1} \sum_{j=1}^n [\hat{f}_\lambda(X_j) - f(X_j)]^2 f(X_j)^{-1} w(X_j).$$

Integrated square error.

$$d_I(\hat{f}_\lambda, f) = \int [\hat{f}_\lambda(x) - f(x)]^2 w(x) dx.$$

Mean integrated square error.

$$d_M(\hat{f}_\lambda, f) = E \int [\hat{f}_\lambda(x) - f(x)]^2 w(x) dx,$$

where $w(x)$ is a nonnegative weight function.

Note that d_M admits the variance-bias square decomposition

$$(4.1) \quad d_M(\hat{f}_\lambda, f) = \int n^{-1} \text{var}[\delta_\lambda(x, X_i)] w(x) dx + \int B(x)^2 w(x) dx,$$

where $B(x)$ denotes the pointwise bias,

$$(4.2) \quad B(x) = \int \delta_\lambda(x, y) f(y) dy - f(x).$$

Marron and Härdle (1986) have shown that, for large n , under reasonable assumptions, these distances are essentially the same in the sense that

$$(4.3) \quad \lim_{n \rightarrow \infty} \sup_{\lambda \in \Lambda_n} \left| \frac{d_A(\hat{f}_\lambda, f) - d_M(\hat{f}_\lambda, f)}{d_M(\hat{f}_\lambda, f)} \right| = 0 \quad \text{a.s.},$$

$$(4.4) \quad \lim_{n \rightarrow \infty} \sup_{\lambda \in \Lambda_n} \left| \frac{d_I(\hat{f}_\lambda, f) - d_M(\hat{f}_\lambda, f)}{d_M(\hat{f}_\lambda, f)} \right| = 0 \quad \text{a.s.},$$

where Λ_n is a finite set whose cardinality grows algebraically fast.

The approximations (4.3) and (4.4) are vital to the theorems of this paper. Other assumptions include the existence of constants $C, C', \delta > 0$ so that

$$(4.5) \quad \#(\Lambda_n) \leq n^C,$$

$$(4.6) \quad C^{-1}n^\delta \leq \lambda \leq Cn^{1-\delta}, \quad \lambda \in \Lambda_n,$$

$$(4.7) \quad w(x) \leq C, \quad x \in \mathbb{R},$$

$$(4.8) \quad f(x) \leq C, \quad x \in \mathcal{S},$$

where \mathcal{S} denotes the support of w ,

$$(4.9) \quad B(x) \leq Cn^{-\delta}, \quad \lambda \in \Lambda_n, \quad x \in \mathcal{S},$$

$$(4.10) \quad \lim_{n \rightarrow \infty} \sup_{\lambda \in \Lambda_n} \left| \frac{\int \text{var}[\sigma_\lambda(x, X_i)] w(x) dx}{C'\lambda} - 1 \right| = 0.$$

Another useful assumption is that for $k = 2, 3, \dots$ there is a constant C_k so that for $m = 2, \dots, k$,

$$(4.11) \quad \int \cdots \int \delta_\lambda(x_{i_1}, x_{j_1}) \cdots \delta_\lambda(x_{i_k}, x_{j_k}) dx_1 \cdots dx_m \leq C_k \lambda^{k-m/2},$$

where $i_1, j_1, \dots, i_k, j_k = 1, \dots, m$ subject to $i_1 \neq j_1, \dots, i_k \neq j_k$, and to each of $1, \dots, m$ appearing at least twice in the list $i_1, j_1, \dots, i_k, j_k$. In the case of kernel estimation (4.11) is a consequence of integration by substitution. Marron and Härdle (1986) show how such a condition is satisfied for the histogram and orthogonal series estimators.

Additional assumptions needed only for the KL cross-validation function include the existence of a constant C so that

$$(4.12) \quad \delta_\lambda(x, x) \leq C\lambda, \quad x \in \mathbb{R},$$

$$(4.13) \quad f(x) \geq C^{-1}, \quad x \in \mathcal{S},$$

$$(4.14) \quad \sup_{j, \lambda, x} |\hat{f}_{\lambda, j^-}(x) - f(x)| \rightarrow 0 \quad \text{a.s.},$$

where $\sup_{j, \lambda, x}$ denotes supremum over $j = 1, \dots, n$, $\lambda \in \Lambda_n$, $x \in \mathcal{S}$.

Observe that (4.13) requires S to be compactly supported. One effect of this is that it avoids the potentially disastrous "tail effects" first reported by Schuster and Gregory (1981). For a detailed analysis of this problem, see Hall (1986a, b). Assumption (4.14) is stated in this form because sufficient conditions for this vary somewhat depending on the particular type of estimator being used. By the calculation (6.1), the uniformity over j may be easily handled. The uniformity over λ and x may be obtained as in Lemma 1 of Härdle and Marron (1985).

Before starting the theorems, it is convenient to define

$$\begin{aligned}
 R &= n^{-1} \sum_{j=1}^n f(X_j)w(X_j) - E[f(X_j)w(X_j)], \\
 (4.15) \quad S &= 2n^{-1} \sum_{j=1}^n [u(X_j)(1 - \log f(X_j)) - R], \\
 T &= - \int f(x)^2 w(x) dx - 2R.
 \end{aligned}$$

It is important to note that R , S and T are independent of λ . For this reason, the fact that both maximizing $\text{KL}(\lambda)$ and minimizing $\text{LS}(\lambda)$ are asymptotically equivalent to minimizing the distances d_A , d_I and d_M is demonstrated by (4.3), (4.4) and

THEOREM 1. *Under the assumptions (4.3), (4.5)–(4.14)*

$$-2n^{-1} \log \text{KL}(\lambda) = d_A(\hat{f}_\lambda, f) + S + o(d_M(\hat{f}_\lambda, f)),$$

in the sense that

$$\lim_{n \rightarrow \infty} \sup_{\lambda \in \Lambda_n} \left| \frac{2n^{-1} \log \text{KL}(\lambda) + d_A(\hat{f}_\lambda, f) + S}{d_M(\hat{f}_\lambda, f)} \right| = 0 \quad a.s.$$

THEOREM 2. *Under the assumptions (4.5)–(4.11)*

$$\text{LS}(\lambda) = d_I(\hat{f}_\lambda, f) + T + o(d_M(\hat{f}_\lambda, f)),$$

in the sense that

$$\lim_{n \rightarrow \infty} \sup_{\lambda \in \Lambda_n} \left| \frac{\text{LS}(\lambda) - d_I(\hat{f}_\lambda, f) - T}{d_M(\hat{f}_\lambda, f)} \right| = 0 \quad a.s.$$

Theorems 1 and 2 are stated in this nonstandard form because this provides the best comparison between KL and LS . Easy consequences of Theorems 1 and 2, respectively, are:

COROLLARY 1. *Under the assumptions (4.3)–(4.14), if $\hat{\lambda}$ is the maximizer of $\text{KL}(\lambda)$ over Λ_n , then*

$$\lim_{n \rightarrow \infty} \frac{d(\hat{f}_{\hat{\lambda}}, f)}{\inf_{\lambda \in \Lambda_n} d(\hat{f}_\lambda, f)} = 1 \quad a.s.,$$

where d is any of d_A , d_I or d_M .

COROLLARY 2. *Under the assumptions (4.3)–(4.11), if $\hat{\lambda}$ is the minimizer of $LS(\lambda)$ over Λ_n , then*

$$\lim_{n \rightarrow \infty} \frac{d(\hat{f}_{\hat{\lambda}}, f)}{\inf_{\lambda \in \Lambda_n} d(\hat{f}_{\lambda}, f)} = 1 \quad a.s.,$$

where d is any of d_A , d_I or d_M .

5. Remarks.

5.1. The main point of this paper is comparison of the KL and LS cross-validation functions. Theorem 1 shows that $KL(\lambda)$ is based on the distance $d_A(\hat{f}_{\lambda}, f)$, while Theorem 2 shows that $LS(\lambda)$ is based on the somewhat more compelling distance $d_I(\hat{f}_{\lambda}, f)$. A more significant advantage of $LS(\lambda)$ is that the term $o(d_M(\hat{f}_{\lambda}, f))$ in Theorem 2 represents error from one source, while in Theorem 1 it represents error from three sources, one of which is the same as that of Theorem 2. A final disadvantage of KL is the stronger assumptions required, especially the uniform convergence assumption (4.14) and the fact that (4.13) requires \mathcal{S} to be compact.

5.2. Despite the poor showing of $KL(\lambda)$ in the above respects, it should be noted that KL and LS are not really comparable because for LS, w must not depend on f , while for KL, $w = u/f$, with u independent of f . This is an advantage of KL because for the important applications of density estimation to discrimination and to minimum Hellinger distance estimation, the latter form is more natural.

5.3. A very important problem in density estimation is that of how to choose between the many available estimators (e.g., histogram, kernel, orthogonal series). Rudemo (1982) has proposed approaching this problem by selecting the estimator which gives the smallest minimum $KL(\lambda)$ [or $LS(\lambda)$]. The second important point of the results of this paper is that they provide theoretical backing to this idea, by indicating that the selected estimator should have the smallest $d_A(\hat{f}_{\lambda}, f)$ [or $d_I(\hat{f}_{\lambda}, f)$, respectively]. Hans-Georg Müller has pointed out that care needs to be taken in interpreting these results with respect to the problem of kernel selection in the kernel estimation case (i.e., use a kernel with several vanishing moments, and hence a higher rate of convergence, or use a nonnegative kernel?). In particular, the $n \rightarrow \infty$ asymptotic results of this paper do not provide proper quantification of the trade off that must occur in kernel selection. It is conjectured that when suitable methods for studying this are found, Rudemo's idea will still be seen to apply.

5.4. As noted in the introduction, the general framework of the results of this paper contain all or part of the results of a number of recent papers as special cases. These include the results of Burman (1985), Hall (1983, 1985), Marron (1985) and Stone (1984, 1985). In most cases the techniques of the present paper

provide a substantial simplification of the proofs in the earlier papers. Also, the unified approach makes it seem easy to provide theoretical backing to some interesting heuristics of Bowman, Hall and Titterton (1984).

5.5. To save space, some of the assumptions of the theorems of this paper have been made more restrictive than necessary. For example, (4.5) can be weakened to Λ_n an interval by a straightforward continuity argument [see Härdle and Marron (1985) for details]. The condition (4.6) can also be substantially weakened [compare Burman (1985), Stone (1984) and Stone (1985)]. Another straightforward extension is to the case of λ vector or matrix valued as discussed by Deheuvels (1977).

6. Proof of Theorems 1 and 2. It is convenient to define, for $j = 1, \dots, n$,

$$\Delta_j = \left[\frac{\hat{f}_{\lambda, j^-}(X_j) - f(X_j)}{f(X_j)} \right] 1_{\mathcal{S}}(X_j), \quad \Delta_j^+ = \left[\frac{f_{\lambda, j^-}^+(X_j) - f(X_j)}{f(X_j)} \right] 1_{\mathcal{S}}(X_j).$$

Note that by (4.13) and (4.14),

$$\sup_{j, \lambda} |\Delta_j^+| \leq \sup_{j, \lambda} |\Delta_j| \rightarrow 0 \quad \text{a.s.},$$

where the suprema are taken over $j = 1, \dots, n$, $\lambda \in \Lambda_n$. For $n = 1, 2, \dots$ define the event

$$U_n = \{ \Delta_j^+ = \Delta_j \text{ for each } \lambda \in \Lambda_n \text{ and each } j = 1, \dots, n \}.$$

Note that

$$\lim_{n \rightarrow \infty} P[U_n] = 1.$$

From the above, it follows that (on the event U_n)

$$\begin{aligned} -2n^{-1} \log \text{KL}(\lambda) - S &= -2n^{-1} \sum_{j=1}^n \left[u(X_j)(1 + \log(1 + \Delta_j)) - \hat{p}(\lambda) - R \right] \\ &= -2n^{-1} \sum_{j=1}^n \left[u(X_j)(1 + \Delta_j) - \hat{p}(\lambda) - R \right] \\ &\quad + d'_A(\hat{f}_\lambda, f) - 2n^{-1} \sum_{j=1}^n r_j u(X_j), \end{aligned}$$

where d'_A is the leave-one-out version of d_A given by

$$d'_A(\hat{f}_\lambda, f) = n^{-1} \sum_{j=1}^n \left[\hat{f}_{\lambda, j^-}(X_j) - f(X_j) \right]^2 f(X_j)^{-1} w(X_j),$$

and where r_j denotes the remainder term of the log Taylor expansion. Theorem 1 follows easily from this, (4.3), (4.14) and the following two lemmas.

LEMMA 1. Under the assumptions of Theorem 1,

$$\lim_{n \rightarrow \infty} \sup_{\lambda \in \Lambda_n} \left| \frac{d_A(\hat{f}_\lambda, f) - d'_A(\hat{f}_\lambda, f)}{d_M(\hat{f}_\lambda, f)} \right| = 0 \quad a.s.$$

LEMMA 2. Under the assumptions (4.5)–(4.11),

$$\lim_{n \rightarrow \infty} \sup_{\lambda \in \Lambda_n} \left| \frac{n^{-1} \sum_{j=1}^n \hat{f}_{\lambda, j-}(X_j)w(X_j) - \int \hat{f}_\lambda(x)f(x)w(x) dx - R}{d_M(\hat{f}_\lambda, f)} \right| = 0 \quad a.s.$$

The proof of Lemma 1 follows in a straightforward manner from

$$\hat{f}_{\lambda, j-}(x) - \hat{f}_\lambda(x) = (n - 1)^{-1} \hat{f}_\lambda(x) - (n - 1)^{-1} \delta_\lambda(x, x)$$

and the assumption (4.12). The proof of Lemma 2 is in Section 7.

An interesting feature of the mathematical structure here is that Lemma 2 contains the hardest part of the proof of Theorem 2 as well. To see this write

$$\begin{aligned} \text{LS}(\lambda) &= d_I(\hat{f}_\lambda, f) - 2n^{-1} \sum_{j=1}^n \hat{f}_{\lambda, j-}(X_j)w(X_j) \\ &\quad + 2n^{-1} \int \hat{f}_\lambda(x)f(x)w(x) dx - \int f(x)^2w(x) dx. \end{aligned} \tag{6.1}$$

Theorem 2 follows easily from this, (4.15) and Lemma 2.

7. Proof of Lemma 2. The conclusion of Lemma 2 may be written as

$$\sup_{\lambda \in \Lambda_n} n^{-1}(n - 1)^{-1} \left| \sum_{i \neq j} U_{i,j} \right| d_M(\hat{f}_\lambda, f)^{-1} \rightarrow 0 \quad a.s.,$$

where

$$\begin{aligned} U_{i,j} &= \delta_\lambda(X_j, X_i)w(X_j) - \int \delta_\lambda(x, X_i)f(x)w(x) dx - f(X_j)w(X_j) \\ &\quad + \int f(x)^2w(x) dx. \end{aligned}$$

For $j = 1, \dots, n$, let

$$W_j = E[U_{i,j}|X_j],$$

and for $i \neq j$ define

$$V_{i,j} = U_{i,j} - W_j.$$

Observe that

$$\begin{aligned} E[V_{i,j}|X_i] &= E[V_{i,j}|X_j] = 0, \\ E[W_j] &= 0. \end{aligned} \tag{7.1}$$

To finish the proof of Lemma 2 it is enough to show that

$$(7.2) \quad \sup_{\lambda \in \Lambda_n} n^{-2} \left| \sum_{i \neq j} V_{i,j} \right| d_M(\hat{f}_\lambda, f)^{-1} \rightarrow 0 \quad \text{a.s.}$$

and that

$$(7.3) \quad \sup_{\lambda \in \Lambda_n} n^{-1} \left| \sum_{j=1}^n W_j \right| d_M(\hat{f}_\lambda, f)^{-1} \rightarrow 0 \quad \text{a.s.}$$

To verify (7.3), note that by the Borel–Cantelli Lemma, it is enough to show that for $\varepsilon > 0$,

$$(7.4) \quad \sum_{n=1}^{\infty} \#(\Lambda_n) \sup_{\lambda \in \Lambda_n} P \left[\left| n^{-1} \sum_{j=1}^n W_j \right| > \varepsilon d_M(\hat{f}_\lambda, f) \right] < \infty.$$

For this, using the notation (4.2), write

$$W_j = B(X_j)w(X_j) - \int B(x)f(x)w(x) dx.$$

From the assumptions (4.7), (4.8) and (4.9) it follows that

$$|W_j| \leq Cn^{-\delta},$$

$$\sigma^2 = \text{var } W_j \leq C^2 \int B(x)^2 w(x) dx.$$

Now Bernstein's inequality [see (2.13) of Hoeffding (1963)] with (in Hoeffding's notation)

$$\lambda = bt/\sigma^2, \quad \tau = nt/b, \quad b = Cn^{-\delta}, \quad t = \varepsilon d_M(\hat{f}_\lambda, f)$$

gives

$$\begin{aligned} P \left[\left| n^{-1} \sum_{j=1}^n W_j \right| > \varepsilon d_M(\hat{f}_\lambda, f) \right] &\leq \exp(-nt^2/2(\sigma^2 + bt/3)) \\ &\leq \exp(-n\varepsilon^2 d_M(\hat{f}_\lambda, f)/2C^2) \\ &\leq \exp(-n^\delta \varepsilon^2/2C^2), \end{aligned}$$

for n sufficiently large. (7.4) is a consequence of this.

To verify (7.2), as in the proof of (7.3) above, together with the Chebyshev inequality, it is enough to show that there is a constant $\gamma > 0$, so that for $k = 1, 2, \dots$ there are constants C_k so that

$$\sup_{\lambda \in \Lambda_n} E \left[n^{-2} \sum_{i \neq j} V_{i,j} d_M(\hat{f}_\lambda, f)^{-1} \right]^{2k} \leq C_k n^{-\gamma k}.$$

But by the cumulant expansion of the $2k$ th centered moment [see, for example,

Kendall and Stuart (1977)], this may be obtained from

$$(7.5) \quad \left| n^{-2k} d_M(\hat{f}_\lambda, \hat{f})^{-k} \sum \text{cum}_k(V_{i_1, j_1}, \dots, V_{i_k, j_k}) \right| \leq C_k n^{-\gamma k},$$

where cum_k is the k th order cumulant and \sum denotes summation over $i_1, j_1, \dots, i_k, j_k = 1, \dots, n$ subject to $i_1 \neq j_1, \dots, i_k \neq j_k$.

To check (7.5), note that by (7.1) and the moment expansion of cum_k , most of the terms in the summation will be 0. In particular, cum_k can be nonzero only when each of $i_1, j_1, \dots, i_k, j_k$ is the same as one of the others. For each such term, let m denote the number of unique elements of $\{1, \dots, n\}$ appearing among $i_1, j_1, \dots, i_k, j_k$. By assumption (4.11) there is a constant C_k so that

$$\left| \text{cum}_k(V_{i_1, j_1}, \dots, V_{i_k, j_k}) \right| \leq C_k \lambda^{k-m/2}.$$

But there is also a constant C_k so that for $m = 2, \dots, k$, the number of nonzero terms in the summation of (7.5) with exactly m distinct indices is bounded by

$$C_k n^m.$$

Hence, by (4.1) and (4.10) there is a constant C_k so that the left side of (7.5) is bounded by

$$C_k n^{-2k} (n^{-1} \lambda)^{-k} \sum_{m=2}^k n^m \lambda^{k-m/2} = C_k \sum_{m=2}^k n^{-k+m} \lambda^{-m/2}.$$

A consequence of this is (7.5). This completes the proof of Lemma 2.

REFERENCES

- BOWMAN, A. W. (1984). An alternative method of cross-validation for the smoothing of density estimates. *Biometrika* **71** 353–360.
- BOWMAN, A. W., HALL, P. and TITTERINGTON, D. M. (1984). Cross-validation in nonparametric estimation of probabilities and probability densities. *Biometrika* **71** 341–352.
- BURMAN, P. (1985). A data dependent approach to density estimation. *Z. Wahrsch. verw. Gebiete* **69** 609–628.
- DEHEUVELS, P. (1977). Estimation non paramétrique de la densité par histogrammes généralisés. *Rev. Statist. Appl.* **25** 5–42.
- FÖLDES, A. and RÉVÉSZ, P. (1974). A general method for density estimation. *Studia Sci. Math. Hungar.* **9** 81–92.
- HABBEMA, J. D. F., HERMANS, J. and VAN DEN BROEK, K. (1974). A stepwise discrimination analysis program using density estimation. In *Compstat 1974: Proceedings in Computational Statistics* (G. Bruckman, ed.) 101–110. Physica Verlag, Vienna.
- HALL, P. (1983). Large sample optimality of least squares cross-validation in density estimation. *Ann. Statist.* **11** 1156–1174.
- HALL, P. (1985). Asymptotic theory of minimum integrated square error for multivariate density estimation. In *Multivariate Analysis VI* (P. R. Krishnaiah, ed.) 289–309. North-Holland, Amsterdam.
- HALL, P. (1986a). On Kullback–Leibler loss and likelihood cross-validation. Unpublished manuscript.
- HALL, P. (1986b). On the estimation of probability densities using compactly supported kernels. Unpublished manuscript.
- HÄRDLE, W. and MARRON, J. S. (1985). Optimal bandwidth selection in nonparametric regression estimation. *Ann. Statist.* **13** 1465–1481.
- HOEFFDING, W. (1963). Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.* **58** 13–30.

- KENDALL, M. G. and STUART, A. (1977). *The Advanced Theory of Statistics: Distribution Theory* 1, 4th ed. Macmillan, New York.
- MARRON, J. S. (1985). An asymptotically efficient solution to the bandwidth problem of kernel density estimation. *Ann. Statist.* **13** 1011–1023.
- MARRON, J. S. and HÄRDLE, W. (1986). Random approximations to some measures of accuracy in nonparametric curve estimation. To appear in *J. Multivariate Anal.*
- RUDEMO, M. (1982). Empirical choice of histogram and kernel density estimators. *Scand. J. Statist.* **9** 65–78.
- SCHUSTER, E. F. and GREGORY, G. G. (1981). On the nonconsistency of maximum likelihood nonparametric density estimators. In *Computer Science and Statistics: Proceedings of the 13th Symposium on the Interface* (W. F. Eddy, ed.) 295–298. Springer, New York.
- STONE, C. J. (1984). An asymptotically optimal window selection rule for kernel density estimates. *Ann. Statist.* **12** 1285–1297.
- STONE, C. J. (1985). An asymptotically optimal histogram selection rule. In *Proc. of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer* (L. M. Le Cam and R. A. Olshen, eds.) 2 513–520. Wadsworth, Monterey, Calif.
- WALTER, G. and BLUM, J. (1979). Probability density estimation using delta sequences. *Ann. Statist.* **7** 328–340.
- WATSON, G. S. and LEADBETTER, M. R. (1965). Hazard analysis II. *Sankhyā Ser. A* **26** 101–116.

DEPARTMENT OF STATISTICS
UNIVERSITY OF NORTH CAROLINA
CHAPEL HILL, NORTH CAROLINA 27514