



The work reported is this document was performed at Lincoln Laboratory, a center for research operated by Massachusetts Institute of Technology, with the support of the Department of the Air Force under Contract F19628-80-C-0002.

This report may be reproduced to satisfy needs of U.S. Government agencies.

The views and conclusions contained in this document are those of the contractor and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the United States Government.

and the general gener

This technical report has been reviewed and is approved for publication.

. . . .

FOR THE COMMANDER

1 P. Luille

Referend L. Loiselle, Lt. Col., USAP Chief, ESD Lincoln Laboratory Project Office

# MASSACHUSETTS INSTITUTE OF TECHNOLOGY LINCOLN LABORATORY

## ON NONPARAMETRIC PROBABILITY DENSITY ESTIMATION USING ORTHOGONAL SERIES

L. K. JONES

Group 92



TECHNICAL NOTE 1980-54

26 NOVEMBER 1980

Approved for public release; distribution unlimited.

LEXINGTON

MASSACHUSETTS

### ABSTRACT

A method of density estimation is proposed, which is a rational modification of orthogonal expansions, combined with a stopping rule determined by a nearest neighbor statistic. This method yields consistent estimates and applies (in principle) to density estimation in any number of dimensions.

and the second se

## ON NONPARAMETRIC PROBABILITY DENSITY ESTIMATION USING ORTHOGONAL SERIES

### I. INTRODUCTION

Among the numerous non-parametric methods of estimating a probability density function, the approximation of this density by a finite fourier series has several computational advantages. Probably the most important of these is the fact that the evaluation of this density at a new data point requires only the storage of certain fourier coefficients. One of the main disadvantages of such an approximation of a density is the difficulty of determining the number of terms in the expansion.

In this note, we propose an approximation which is a rational function of a finite fourier series. The number of terms in this series depends, in a very natural way, on the nearest neighbor error rate for the sample data when compared to a sample drawn from a reference distribution. In III we show that the method is consistent and in IV we remark on the relevance of this method in hypothesis testing.

II. SECOND ORDER SOLUTION TO THE BINARY DECISION PROBLEM

Let  $p_1$ ,  $p_2$  be two Lebesgue measurable, bounded ( $\leq K$ ) probability density functions on the unit cube, I, in  $\mathbb{R}^n$ . We assume further that  $p_1 \neq p_2$  on some set of positive measure in I and that for some  $\delta \geq 0$ ,  $p_1 \geq \delta$  on I. Let  $\mathbf{L} = \{f \in L_2(I) : E_1 f = \int f p_1 dx = 0, E_2 f = \int f p_2 dx = 1\}$ . According to [1], a second order solution for an optimal discriminant  $\overline{f} \in \mathbf{L}$ , for the binary hypothesis test  $H_1$ : X has density  $p_1$  vs.  $H_2$ : X has density  $p_2$ , is a critical point for some real  $\alpha$  of the functional

$$J_{\alpha}(f) = \alpha VAR_{1} f + (1-\alpha) VAR_{2} f$$
 (1)

In fact if we restrict ourselves to the case  $0 < \alpha < 1$  and solve (1) for the unique (to within a null function) critical point (and global minimum of  $J_{\alpha}(f)$  for  $f \in \mathcal{L}$ ), we obtain by elementary variational calculus

$$\overline{\mathbf{f}} = \frac{\left[ (1-\alpha) - \lambda \right] \mathbf{p}_2 / \mathbf{p}_1 + \lambda}{\alpha + (1-\alpha) \mathbf{p}_2 / \mathbf{p}_1}$$
(2)

with

$$0 > \lambda = \frac{(1-\alpha)\int \frac{p_2 p_1}{\alpha p_1 + (1-\alpha) p_2} dx}{\int \frac{(p_2 - p_1) p_1}{\alpha p_1 + (1-\alpha) p_2} dx} = (1-\alpha) - J_{\alpha}(\overline{f})$$

It follows that  $\overline{f}$  is rational and increasing in  $(p_2/p_1)$ and hence optimal (by an adjustment of threshold) for minimum total error (or Neyman-Pearson at level  $\beta$ ) hypothesis testing.

#### III. DENSITY ESTIMATION

For simplicity we consider only the case  $\alpha = \frac{1}{2}$ . Similar results may be obtained for other  $\alpha$ . Solving (2) for  $p_2/p_1$  we obtain

$$\frac{\mathbf{p}_2}{\mathbf{p}_1} = \frac{\frac{1}{2}\overline{\mathbf{f}} - \lambda}{(\frac{1}{2} - \lambda) - \frac{1}{2}\overline{\mathbf{f}}} = \frac{\overline{\mathbf{f}} - 1 + 2\mathbf{J}_{\frac{1}{2}}(\overline{\mathbf{f}})}{2\mathbf{J}_{\frac{1}{2}}(\overline{\mathbf{f}}) - \overline{\mathbf{f}}}$$
(3)

We now write

$$J_{\frac{1}{2}}(\overline{f}) = \frac{1}{2} + \frac{\varepsilon_{nn}}{4(\frac{1}{2}-\varepsilon_{nn})}$$
(4)

where  $\varepsilon_{nn} = \int \frac{p_2 p_1}{p_1 + p_2} dx$  is known as the limiting nearest neighbor error rate, i.e., if we generate n independent class 1 samples from a distribution with density  $p_1$  and similarly n class 2 samples from a distribution with density  $p_2$  and then classify new samples (drawn from class 1 or class 2 with equal probability) as the class of the (a) nearest neighbor in the original 2n, then the classification error of this procedure approaches  $\varepsilon_{nn}$ as  $n \rightarrow \infty$  with probability 1. (See [2].)

We now make a final assumption that  $p_1 \equiv 1$  on I. Again, results analagous to the following will still hold provided  $p_1$ is strictly bounded away from 0 in I.

Suppose we are given n independent samples from a distribution with density  $p_2$ . Let  $l \equiv \varphi_1$ ,  $\varphi_2$ ,... be a complete orthonormal system for  $L_2(I)$ . Finally, let  $v_n$  be the empirical density determined by the n sample points. Now, consider the solution

of the variational problem: minimize

$$\frac{1}{2} \text{VAR}_{1} f + \frac{1}{2} \text{VAR}_{v_{n}} f = J_{N}^{n} (f)$$
 (\*)

such that

$$f = \sum_{i=1}^{N} a_{i} \varphi_{i}$$
$$E_{1} f = 0$$
$$E_{v_{n}} f = 1$$

where N is determined by a "stopping rule". We then let  $f_n$  be the above minimizing f. Before describing the determination of N, we show that the preceding variational problem has solutions with probability 1 for large enough N.

<u>Lemma</u> Assume n is fixed. Then with probability one (\*) has solutions for large enough N. In fact min  $J_N^n(f) \to 0$  as  $N \to \infty$  with probability one.

<u>Proof:</u> Let L be any positive integer and  $\varepsilon = 0$ . Then there is an N<sub>0</sub> such that, for N>N<sub>0</sub>, there are functions  $\Psi_1, \Psi_2, \dots, \Psi_L$  $\varepsilon \langle \Psi_1, \Psi_2, \dots, \Psi_N \rangle$  with the properties:

- (i)  $||\Psi_{i}||_{2}^{2} \leq \frac{1}{L} + \epsilon$
- (ii) there exist disjoint subsets  $A_1, \dots A_2$  with  $m(\bigcup A_i) > 1-\varepsilon$  s.t.  $x \in A_i$  implies 1

$$|\Psi_{i}(\mathbf{x})-1| \leq \varepsilon$$
 and  $|\Psi_{i}(\mathbf{x})| \leq \varepsilon$  ( $j \neq i$ ).

Hence, with probability (wrt  $p_2$ ) > (1-K $\epsilon$ )<sup>n</sup>, each of our samples  $\ell$  will lie in some  $A_{i_{\ell}}$ . Let us now consider the function

$$\widetilde{\mathbf{f}} = \frac{\sum_{\ell=1}^{n} \Psi_{i_{\ell}} - \sum_{\ell=1}^{n} \int \Psi_{i_{\ell}} \, \mathrm{d}\mathbf{x}}{\frac{1}{n} \sum_{k=1}^{n} \sum_{\ell=1}^{n} \Psi_{i_{\ell}}(\mathbf{x}_{k}) - \sum_{\ell=1}^{n} \int \Psi_{i_{\ell}} \, \mathrm{d}\mathbf{x}}$$

Clearly  $E_1 \quad \tilde{f}=0$ ,  $E_{v_n} \quad \tilde{f}=1$ . We have further

$$\operatorname{VAR}_{1} \widetilde{f} = ||\widetilde{f}||_{2}^{2} \leq \left(\frac{n\sqrt{\frac{1}{L}+\varepsilon}}{1-n\varepsilon-n\sqrt{\frac{1}{L}+\varepsilon}}\right)^{2}$$

$$\operatorname{VAR}_{\nu_{n}} \widetilde{f} \leq \left(\frac{2n\varepsilon}{1-n\varepsilon-n\sqrt{\frac{1}{L}}+\varepsilon}\right)^{2}$$

Hence,  $J_N^{n}(\tilde{f})$  becomes arbitrarily small as  $N \rightarrow \infty$  with probability arbitrarily close to one.

The solutions of (\*) can be easily obtained by the method of Lagrange multipliers. Since  $\varphi_1 \equiv 1, (*)$  is reduced to solving the following for  $a_i = 1$ 

$$\min \begin{bmatrix} \frac{1}{2} \sum_{i=1}^{N} a_{i}^{2} + \frac{1}{2} \sum_{i,j=2}^{N} a_{i}a_{j} \overline{\varphi}_{ij} \end{bmatrix}$$

such that

$$\sum_{2}^{N} a_{i} \overline{\varphi}_{i} = 1$$

$$\overline{\varphi}_{i} = \frac{1}{n} \sum_{1}^{n} \varphi_{i}(x_{\ell})$$

$$\overline{\varphi}_{ij} = \frac{1}{n} \sum_{1}^{n} \varphi_{i}(x_{\ell}) - \varphi_{j}(x_{\ell})$$

For the determination of N=N<sub>p</sub>, we first estimate  $J_{k}(\overline{f})$  by

$$\overline{J}_{n} = \frac{1}{2} + \frac{\varepsilon^{n}}{4(\frac{1}{2}-\varepsilon^{n})}$$
(5)

where  $\varepsilon^n$  is the expected nearest neighbor error rate of the n samples with the leaving-one-out method:

$$\epsilon^{n} = \frac{1}{n} \sum_{\ell=1}^{n} \left[ 1 - (1 - V_{\ell})^{n-1} \right]$$
(6)

where  $V_{\ell}$  is the volume of the intersection of I with a sphere centered at  $x_{\ell}$  and of radius equal to the distance between  $x_{\ell}$ and its nearest neighbor in  $\{x_k\}_{k\neq\ell}$ . Now  $\varepsilon^n \to \varepsilon_{nn}$  with probability one as  $n \to \infty$  and hence  $\overline{J}_n \to J_{\frac{1}{2}}(\overline{f})$  with probability one as  $n \to \infty$ . Let  $N_n$  be an N which minimizes  $|J_N^{\ n}(f_n) - \overline{J}_n|$ . By the lemma such an N exists with probability one provided that  $\overline{J}_n > 0$  and this is true with probability one as  $n \to \infty$ .

The estimate we then use for  $p_2$  is

$$\mathbf{p}_{n} = \frac{\mathbf{f}_{n} - \mathbf{1} + 2\overline{\mathbf{J}}_{n}}{2\overline{\mathbf{J}}_{n} - \mathbf{f}_{n}} \qquad . \tag{7}$$

If we should know the value of K, we may use the truncated estimate

$$\hat{\mathbf{p}}_{n} = (\mathbf{p}_{n} \vee \mathbf{0}) \wedge \mathbf{K}$$
(8)

We now make the following consistency claim.

<u>Theorem</u>  $p_n \rightarrow p_2$  in Lebesgue measure with probability one and  $\hat{p}_n \xrightarrow{L_2} p_2$  with probability one.

<u>Proof</u> From the form of (3), (7), (8) and the fact that  $\overline{J}_n \rightarrow J_{\frac{1}{2}}(\overline{f})$ , it suffices to show that  $f_n \xrightarrow{L_2} \overline{f}$  with probability one.

Note that  $\varphi_2$ ,  $\varphi_3$ ,... are linearly independent and dense in  $\{f: \int fdx=0\} \cap L_2(\frac{1}{2}+\frac{p_2}{2})$  where  $L_2(\frac{1}{2}+\frac{p_2}{2})$  denotes the set of square integrable functions wrt. to a measure whose density is  $\frac{1}{2}+p_2/2$ . Now form a complete orthonormal basis  $\xi_2$ ,  $\xi_3$ ,... of  $\{f: \int fdx=0\} \cap L_2(\frac{1}{2}+\frac{p_2}{2})$  where each  $\xi_i$  is a linear combination of  $\varphi_2$ ,  $\varphi_3$ ,...  $\varphi_i$ . Let  $c_i = \int \xi_i p_2$ . Then  $\overline{f} = \sum_2^{\infty} b_i \xi_i$  where  $b_i$  is the solution of min  $\sum_2^{\infty} b_i^2$  such that  $\sum_2^{\infty} c_i b_i = 1$ . This is just  $b_i = c_i / \sum_2^{\infty} c_i^2$ .

Similarly, we form a complete orthonormal basis  $n_2^n$ ,  $n_3^n$ ,... of  $\{f: \int f=0\} \cap L_2(\frac{1}{2} + \frac{\nu_n}{2})$ , with each  $n_i^n$  a linear combination of  $\varphi_2, \varphi_3, \dots \varphi_i$ . Let  $d_i^n = \int n_i^n \nu_n$ . The solution of (\*) is given by  $f_n = \sum_{2}^{N_n} d_i^n n_i^n / \sum_{2}^{N_n} (d_i^n)^2$ .

Clearly,  $d_i^n \to c_i$  and  $|| n_i^n - \xi_i ||_2 \to 0$  with probability one as  $n \to \infty$  for each i. Since

$$\left| \left( \sum_{2}^{N} c_{1}^{2} \right)^{-1} - \left( \sum_{2}^{N} (d_{1}^{n})^{2} \right)^{-1} \right| \to 0$$

with probability one as  $n\!\rightarrow\!\infty$  , for each N, it follows that

$$\left| \left( \sum_{2}^{N_{n}} (\mathbf{a}_{1}^{n})^{2} \right)^{-1} - \left( \sum_{2}^{\infty} \mathbf{c}_{1}^{2} \right)^{-1} \right| \to 0$$

with probability one as  $n \to \infty$ , and hence  $\sum_{i=1}^{N} (d_i^{n-1})^2 \to \sum_{i=1}^{\infty} c_i^{n-2}$ with probability one as  $n \to \infty$ . Finally, for any  $\varepsilon > 0$  pick M

such that  $\sum_{M}^{\infty} c_i^2 < \epsilon$ . Then

$$\overline{\lim} || f_n - \overline{f} ||_2 \leq \overline{\lim} || \sum_{M}^{N_n} d_i^n \eta_i^n ||_2 + || \sum_{M}^{\infty} c_i \xi_i ||_2$$

with probability one. But

$$\|g\|_{2} \leq \sqrt{2} \|g\|$$
 and  $\|g\|_{2} \leq \sqrt{2} \|g\|_{1}^{2}$   
 $\frac{1}{2} + \frac{p_{2}}{2}$   $\frac{1}{2} + \frac{n}{2}$ 

Hence  $\overline{\lim} \|f_n - \overline{f}\|_2 \le 2\sqrt{2\varepsilon}$  with probability one. Since  $\varepsilon$  was arbitrary, the proof is complete.

### IV. A REMARK ON HYPOTHESIS TESTING

If we are given two sets of data  $\{x_i\}_{i=1}^n$  and  $\{y_j\}_{j=1}^n$  and are then given the task of constructing an optimal discriminant between the two classes, we might solve

$$\min \frac{1}{2} \begin{bmatrix} VA_1 + \frac{1}{2} VAR_{v_n} \end{bmatrix}$$
 (\*\*)

such that

$$E_{\mu_{n}} f=0$$

$$E_{\nu_{n}} f=1$$

$$f = \sum_{i}^{N_{n}} \varphi_{i}$$

$$\mu_{n}, \nu_{n} \text{ empirical densities for}\{x_{i}\}, \{y_{j}\}$$

where N<sub>n</sub> is determined by the analagous "stopping rule".

Unfortunately the consistency proof does not apply in this case since we are unable to find a bound for  $||g||_{p_1+p_2}$  in terms of  $||g||_{\frac{v_n+\mu_n}{2}}$ . Hence, the reference density  $p_1 \equiv 1$  "forced" the consistency. We therefore recommend the estimates:  $\hat{p}_1$  by the method of III using  $\{x_i\}$  and then similarly  $\hat{p}_2$  using  $\{y_j\}$ . The discriminant  $\hat{p}_2/\hat{p}_1$  will then be optimal with probability one

as n→∞.

#### REFERENCES

- [1]. L. Jones, "K<sup>th</sup> Order Solutions to the Problem of Finding Optimal Discriminant Functions", Technical Note 1980-4, Lincoln Laboratory, M.I.T. (10 January 1980), submitted to S.I.A.M. J. Appl. Math. DTIC AD-A0822396/3.
- [2]. T. Cover and P. Hart, "Nearest Neighbor Pattern Classification", IEEE Trans. Inf. Theory <u>IT-13</u>, pgs. 21-27 (1967).

19 REPORT DOCUMENTATION PAGE	READ INSTRUCTIONS BEFORE COMPLETING FORM
REPORT NUMBER 2. GOVT ACCES	SION NO. 3. RECIPIENT'S CATALOG NUMBER
ESD TR-89-230 / AD AC94	728
TITLE (and Subsiste)	5. TYPE OF REPORT & PERIOD COVERED
On Nonparametric Probability Density Estimation	( ] Technical Note
Using Orthogonal Series	6. PERFORMING ORG. REPORT NUMBER
	Technical Note 1980-54
Lee K. Jones	F19628-80-C-0002
PERFORMING ORGANIZATION NAME AND ADDRESS	10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
Lincoln Laboratory, M.I.T. P.O. Box 73 Lexington, MA 02173 (6) (27A	Program Element No. 63311F Project No. 627A
1. CONTROLLING OFFICE NAME AND ADDRESS	12. REPORT DATE
Air Force Systems Command, USAF Andrews AFB	( / /. 26 Nov <del>ember 19</del> 80
Washington, DC 20331	TJ. NUMBER OF PAGES
4. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)	15. SECURITY CLASS. (of this report)
Electronic Systems Division	Unclassified
Bedford, MA 01731	15a. DECLASSIFICATION DOWNGRADING SCHEDULE
Approved for public release; distribution unlimited 7. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different fre	om Report)
Approved for public release; distribution unlimited 7. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different fro 8. SUPPLEMENTARY NOTES	om Report;
Approved for public release; distribution unlimited 7. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different fro 8. SUPPLEMENTARY NOTES None	om Report)
Approved for public release; distribution unlimited 7. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different fr 8. SUPPLEMENTARY NOTES None 9. KEY WORDS (Continue on reverse side if necessary and identify by block number	om Report)
Approved for public release; distribution unlimited 7. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different fro 8. SUPPLEMENTARY NOTES None 9. KEY WORDS (Continue on reverse side if necessary and identify by block number density estimation Fourier coefficient	om Report) ') nts orthogonal expansions
Approved for public release; distribution unlimited 7. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different fri 8. SUPPLEMENTARY NOTES None 9. KEY WORDS (Continue on reverse side if necessary and identify by block number density estimation Fourier coefficien 9. ABSTRACT (Continue on reverse side if necessary and identify by block number	om Report) ') nts orthogonal expansions
Approved for public release; distribution unlimited 7. DISTRIBUTION STATEMENT (of the obstract entered in Block 20, if different for 8. SUPPLEMENTARY NOTES None 9. KEY WORDS (Continue on reverse side if necessary and identify by block number density estimation Fourier coefficien 0. ABSTRACT (Continue on reverse side if necessary and identify by block number, A method of density estimation is proposed, which expansions, combined with a stopping rule determined method yields consistent estimates and applies (in prin ber of dimensions,	on Repon; "," nts orthogonal expansions ," h is a rational modification of orthogonal by a mearest neighbor statistic. This nciple) to density estimation in any num-
Approved for public release; distribution unlimited . DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different for . SUPPLEMENTARY NOTES None . KEY WORDS (Continue on reverse side if necessary and identify by block number density estimation Fourier coefficien . ABSTRACT (Continue on reverse side if necessary and identify by block number A method of density estimation is proposed, which expansions, combined with a stopping rule determined method yields consistent estimates and applies (in prin ber of dimensions. FORM 1473 EDITION OF 1 NOV 65 15 OBSOLETE	on Repon)  ' ' nts orthogonal expansions  ' ' A is a rational modification of orthogonal by a nearest neighbor statistic. This nciple) to density estimation in any num- UNCLASSIFIED

A REAL PROPERTY AND AND

4

5.0

8 · . . .

