

AD-A104 860

RICE UNIV HOUSTON TEX DEPT OF MATHEMATICAL SCIENCES
AN ALGORITHM FOR NONPARAMETRIC DENSITY ESTIMATION; (U)
MAY 76 D W SCOTT, R A TAPIA, J R THOMPSON

F/G 12/1

E-(40-1)-5046

UNCLASSIFIED

NL

1 of 1
AD-A104 860



END
DATE
FILMED
10 14
DTIC

LEVEL

For: Computer Science and Statistics: NINTH ANNUAL Symposium on the Interface.

May 1976

DTIC ELECTE S OCT 1 1981 D

2

ADA104800

AN ALGORITHM FOR NONPARAMETRIC DENSITY ESTIMATION

David W. Scott Dept. of Biostatistics Baylor College of Medicine Houston, Texas 77030

Richard A. Tapia Dept. of Mathematical Sciences Rice University Houston, Texas 77001

James R. Thompson Dept. of Mathematical Sciences Rice University Houston, Texas 77001

ABSTRACT

A numerical algorithm is given for implementing a nonparametric maximum penalized likelihood estimator similar to those proposed by Good and Gaskins and those proposed by de Montricher, Tapia and Thompson. It is shown how the resulting nonlinear constrained optimization problem may be effectively solved by using Tapia's approach to Newton's method for constrained problems.

1. Introduction. de Montricher, Tapia and Thompson demonstrated that the standard histogram was an unstable maximum likelihood density estimator and considered maximum penalized likelihood estimators similar to those previously considered by Good and Gaskins (1971). Specifically suppose we are given the random sample x1, ..., xn in (a, b). Let H0(a, b) consist of the functions f defined on (a, b) with the property that f(a) = f(b) = 0 and f' is a member of L2(a, b). Estimate the density function which gave rise to the random sample x1, ..., xn by the solution of the constrained optimization problem

(1.1) max L(f); f in H0(a, b), f >= 0 and

int_a^b f(x) dx = 1,

where

(1.2) L(f) = int_a^b f(x) exp(-alpha int_a^b |f'(x)|^2 dx), (alpha > 0).

The functional L in (1.2) is called the penalized likelihood and the solution of (1.1) is called the maximum penalized likelihood estimator based on the random sample x1, ..., xn. de Montricher, Tapia and Thompson (1975) proved that problem (1.1) has a unique solution and is a monospline of degree two,

i.e., a polynomial of degree two plus a spline of degree one. We now give a numerical algorithm for approximating this monospline.

2. The Discrete Problem. For given n, consider the mesh t0, ..., tn+1 where ti = a + ih, i = 0, ..., n+1 with h = (b-a)/(n+1). Let H1 denote the vector space of all continuous piecewise linear functions which have knots at t1, ..., tn and vanish at a and b. For p in H1 let yi = p(ti), i = 0, ..., n+1. Then y0 = yn+1 = 0 and

(2.1) p(x) >= 0 <=> yi >= 0, i = 1, ..., n

(2.2) int_a^b p(x) dx = h sum_{i=0}^n yi

(2.3) int_a^b p'(x)^2 dx = 1/h sum_{i=0}^n (yi+1 - yi)^2

Let

(2.4) v1 = # of xi in [a, t1 + h/2)

(2.5) vi = # of xi in [ti-1 + h/2, ti + h/2), i = 2, ..., n-1

(2.6) vn = # of xi in [tn-1 + h/2, b).

We shall assume that we have enough data so that vi > 0 for all i. Our finite dimensional

This work was supported in part by ONR grant ONR-042-283 and ERDA contract E-(48-1)-5646

81 9 30 078

This document has been approved for public release and sale; its distribution is unlimited.

DTIC FILE COPY

15

4-573

estimator based on a sample of size 100. In Figures 3 and 4 we show the D.M.P.L.E. estimators for the 50-50 mixture of two normal distributions, both having variance 1 and with means at -1.5 and +1.5 for samples of size 25 and 100 respectively.

One comforting feature of the maximum penalized likelihood procedure is the relatively robust quality of the estimator in that changes of the optimal α with N and from distribution to distribution tend not to be traumatic, and that a rough and ready guess for α (e.g., 10) is frequently satisfactory. In Figures 5 and 6 we show an estimate for the Gaussian mixture mentioned above for a sample size of 300 and α values of 10 and .1.

If the density to be estimated is denoted by $f(\cdot)$ and the D.M.P.L.E. is denoted by $\hat{f}(\cdot)$, then we consider as one measure of estimate quality the average integrated mean square error

$$(4.3) \text{IMSE} = \int (\hat{f}(x) - f(x))^2 f(x) dx .$$

I.M.S.E.'s for various α and N are given in Table 1 for the standard normal, the 50-50 normal mixture mentioned above, the t distribution with 5 degrees of freedom and the F distribution with (10,10) degrees of freedom.

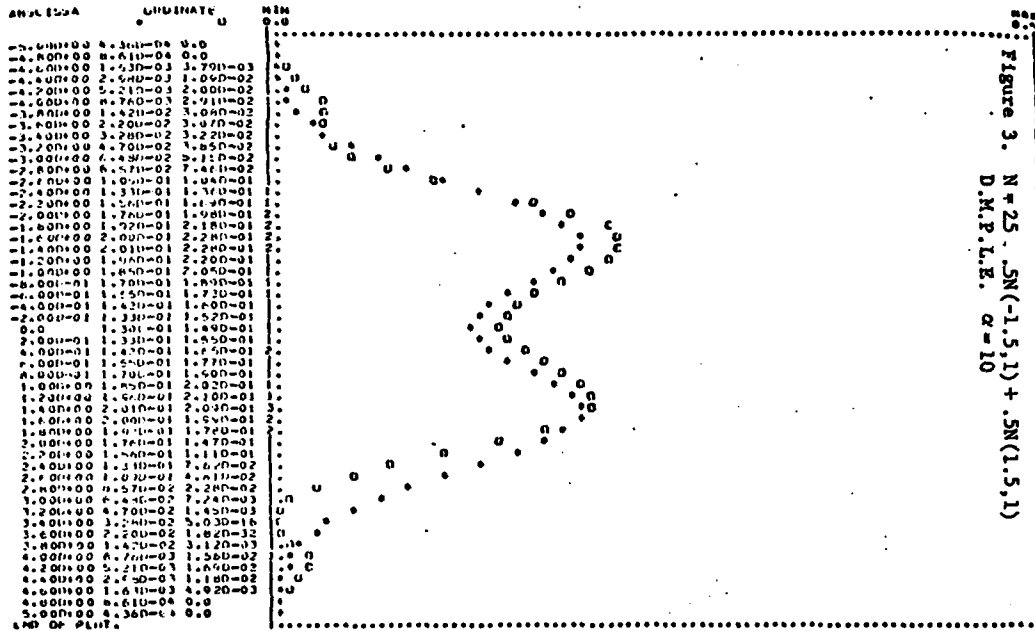
REFERENCES

1. Good, I.J. and Gaskins, R.A. (1971). "Nonparametric roughness penalties for probability densities," Biometrika 58, pp. 255-277.
2. de Montricher, G., Tapia, R.A., and Thompson, J.R. (1975). "Nonparametric maximum likelihood estimation of probability densities by penalty function methods," Annals of Statistics 3, pp. 1329-1348.
3. Tapia, R.A. (1974). "A stable approach to Newton's method for general mathematical programming problems in R^n ," Journal of Optimization Theory and Applications 14, pp. 453-476.
4. Tapia, R.A. (1976). "Diagonalized multiplier methods and quasi-Newton methods for constrained optimization," to appear in Journal of Optimization Theory and Applications.

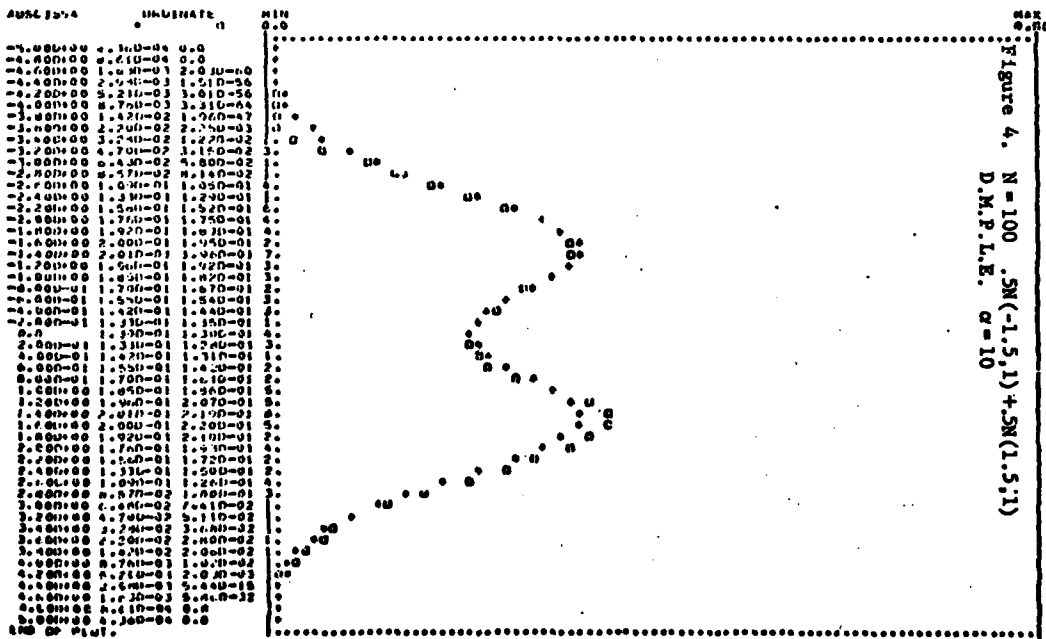
TABLE 1
Average I.M.S.E. of the D.M.P.L.E. for α Perturbed by a Factor of Two. Divide α by 10 for the $F_{10,10}$ Samples.

Sample	I.M.S.E. for		
	$\alpha = 5$	$\alpha = 10$	$\alpha = 20$
$N(0.1) N = 25$.00242	.00267	.00427
$N(0.1) N = 100$.00093	.00079	.00089
$N(0.1) N = 400$.00037	.00033	.00035
$N(0.1) N = 800$.00028	.00022	.00019
Bimodal $N = 25$.00197	.00159	.00152
Bimodal $N = 100$.00070	.00054	.00171
Bimodal $N = 400$.00030	.00024	.00022
$t_5 N = 25$.00297	.00282	.00350
$t_5 N = 100$.00092	.00084	.00101
$t_5 N = 400$.00039	.00032	.00030
$F_{10,10} N = 25$.03208	.03865	.05519
$F_{10,10} N = 100$.00996	.01390	.02411
$F_{10,10} N = 400$.00292	.00450	.00740

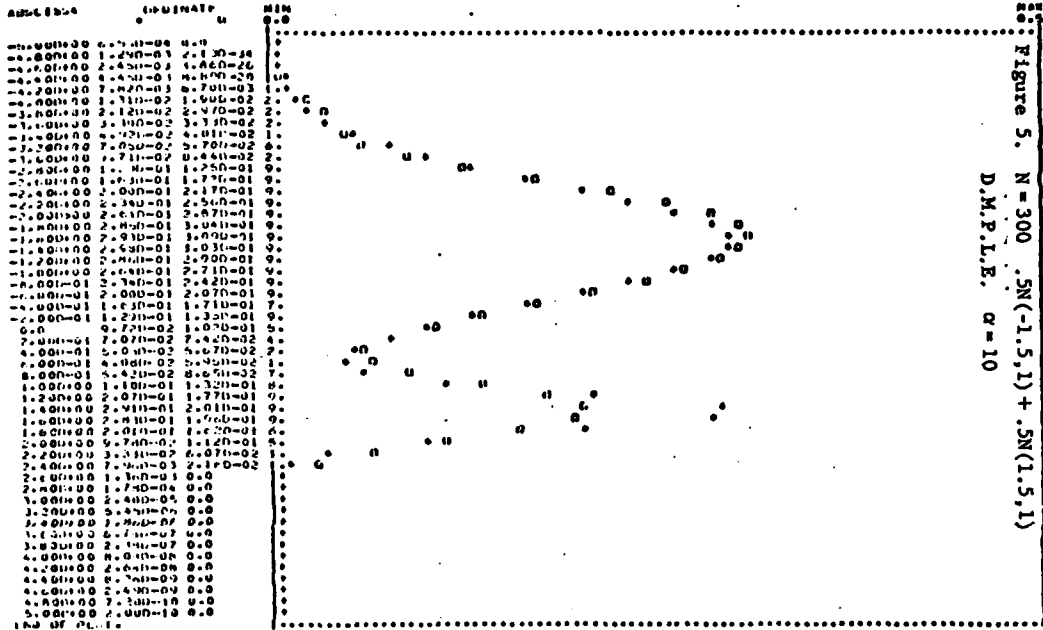
MINIMUM SQUARE WITH MEANS = 1.50 VARIANCE OF LEFT = 1 WITH WEIGHT AND VARIANCE OF RIGHT = 0.5000 1.0000
 SAMPLE SIZE = 25 WITH 50 SAMPLES ON RIGHT VS.
 UNBIASED MAXIMUM LIKELIHOOD ESTIMATE
 WITH WEIGHTING PARAMETER ALPHA = 0.10000000000002
 41 MESH POINTS FROM -5.0000 TO 5.0000
 MESH SIZE INTERVAL 0.25000
 INTEGRATED MEAN SQUARE ERROR 0.7011974300-03
 INTEGRATED SQUARE ERROR 0.6201120600-02
 MAXIMUM ABSOLUTE DIFFERENCE 0.3477620520-01
 LOG LIKELIHOOD TERM -0.4521376360-02
 LOG PENALTY TERM -0.23516194510-01



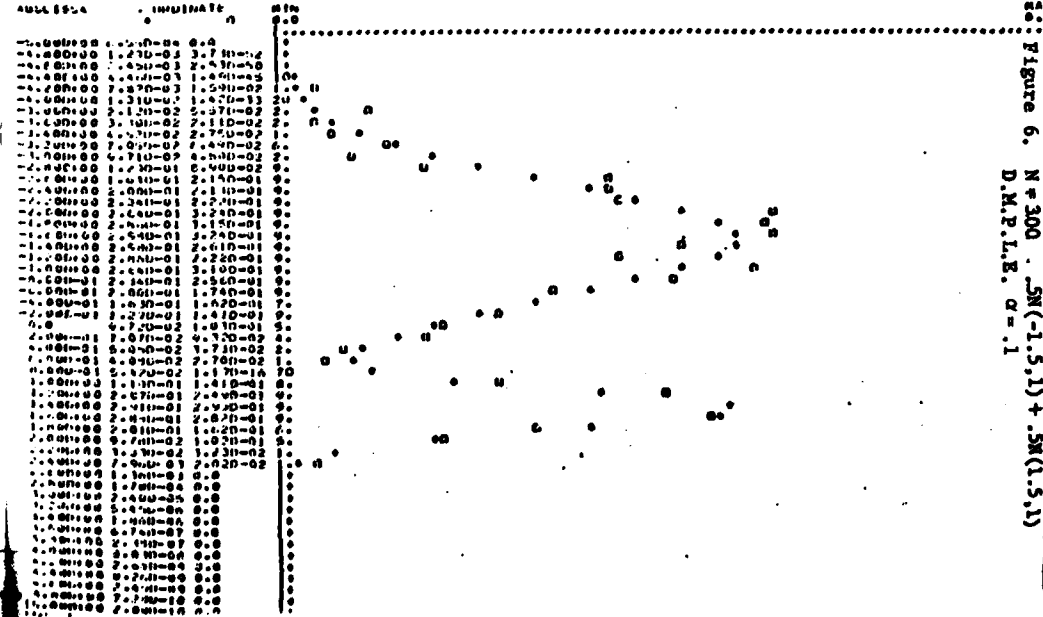
MINIMUM SQUARE WITH MEANS = 1.50 VARIANCE OF LEFT = 1 WITH WEIGHT AND VARIANCE OF RIGHT = 0.5000 1.0000
 SAMPLE SIZE = 100 WITH 50 SAMPLES ON RIGHT VS.
 UNBIASED MAXIMUM LIKELIHOOD ESTIMATE
 WITH WEIGHTING PARAMETER ALPHA = 0.10000000000002
 41 MESH POINTS FROM -5.0000 TO 5.0000
 MESH SIZE INTERVAL 0.25000
 INTEGRATED MEAN SQUARE ERROR 0.1321896400-01
 INTEGRATED SQUARE ERROR 0.5276400620-01
 MAXIMUM ABSOLUTE DIFFERENCE 0.2323925080-01
 LOG LIKELIHOOD TERM -0.1867027490-03
 LOG PENALTY TERM -0.17131662430-01



MINIMUM MAXIMUM WITH MEANS $\alpha = 0.1$ VARIANCE OF LEFT = 1 WITH WEIGHT AND VARIANCE OF RIGHT = 0.2500 0.1111
 SAMPLE SIZE = 100 WITH 75 SAMPLES ON EACH VS.
 DISTRICTION MAXIMUM LEFT INHOD PENALIZED ESTIMATE
 ESTIMATED MAXIMUM LEFT INHOD PENALIZED ESTIMATE
 41 MEAN POINTS FROM = 20000 TO 20000
 DISTRICTION MAXIMUM LEFT INHOD PENALIZED ESTIMATE
 ESTIMATED MAXIMUM LEFT INHOD PENALIZED ESTIMATE
 MAXIMUM ABSOLUTE DIFFERENCE
 LOG PPMAL TV FROM



MINIMUM MAXIMUM WITH MEANS $\alpha = 0.1$ VARIANCE OF LEFT = 1 WITH WEIGHT AND VARIANCE OF RIGHT = 0.2500 0.1111
 SAMPLE SIZE = 100 WITH 75 SAMPLES ON EACH VS.
 DISTRICTION MAXIMUM LEFT INHOD PENALIZED ESTIMATE
 ESTIMATED MAXIMUM LEFT INHOD PENALIZED ESTIMATE
 41 MEAN POINTS FROM = 20000 TO 20000
 DISTRICTION MAXIMUM LEFT INHOD PENALIZED ESTIMATE
 ESTIMATED MAXIMUM LEFT INHOD PENALIZED ESTIMATE
 MAXIMUM ABSOLUTE DIFFERENCE
 LOG PPMAL TV FROM



END

DATE
FILMED

10-81

DTIC