------------------------
DEPARTMENT OF STATISTICS
------------------------
University of Wisconsin
1210 W. Dayton St.
Madison, WI  53706

TECHNICAL REPORT NO. 711

June 1983

# CROSS VALIDATED SPLINE METHODS FOR THE ESTIMATION OF THREE DIMENSIONAL TUMOR SIZE DISTRIBUTIONS FROM OBSERVATIONS ON TWO DIMENSIONAL CROSS SECTIONS

by

Douglas Nychka
and
Grace Wahba

Department of Statistics

and

Stanley Goldfarb
and
Thomas Pugh

Department of Pathology and Laboratory Medicine

# ABSTRACT

We study the problem of estimating the distribution of the three dimensional radii of a collection of spheres, given measurements of the two dimensional radii of a sample of planar cross sections. This problem arises in the estimation of the tumor size distribution of spherical microtumors induced in mouse livers following injection of a carcinogen. We first convert this problem to a form suitable for the application of cross validated spline methods for the solution of ill posed integral equations given noisy data. Then we develop special numerical techniques which will allow the spline methods to be accurately applied to integral equations like those associated with the present problem. We apply the resulting method to some mouse liver data. The subject mouse liver has been completely dissected, allowing a rare comparison of the estimate with the "truth". The statistical properties of the estimate are explored via Monte Carlo methods. The interplay between statistical and numerical analytic methods for problems like this are explored and the use of eigensequence plots for studying ill posedness is described.

Key words:    random spheres model
              stereology
              cross validated splines
              tumor size distribution
              ill posed problem

## 1. INTRODUCTION

We have been working with data from experiments in pathology studying the growth of micro-tumors (hepatocellular foci) in the livers of mice (see Koen, H., Pugh, T. and Goldfarb, S. (1983)). Mice are injected at 15 days of age with a carcinogen which induces the formation of malignant tumors in the liver. After a fixed period of time the mice are sacrificed and samples of liver tissue stained and embedded in a paraffin block. The matrix of paraffin enables the sample to be sliced very thinly, and these slices are mounted on microscope slides. Tumors in the sample will now appear in cross section on these slides and their cross sectional area or radii, if spherical, can be measured.

It is desired to estimate the number density and three dimensional size distribution of the liver tumors from the cross sectional observations. In these particular experiments, a single mouse liver may contain anywhere from a few to several hundred micro-tumors. Different mathematical models for tumor growth have different implications for the variation of tumor size distribution with mouse age. Thus, it is desired to identify tumor size distributions for groups of experimental animals sacrificed at different times after the exposure to the carcinogen. These growth models are important because they might suggest some of the mechanisms which initiate and promote liver cancer. By the limitations of the dissection procedure, tumors can only be identified by their cross sections. Since tumors of different sizes can produce the same size cross sections, there is not a direct correspondence between the cross sectional data and the distribution

of tumor sizes. Although it is possible to take many, closely spaced slices and completely reconstruct each tumor, this procedure is both tedious and costly. What is required is a statistical method that estimates the 3 dimensional tumor size distribution from observations of two dimensional cross sections from a modest number of slices.

The biology of the liver suggests that the tumors will be uniformly distributed throughout the tissue, while examination of successive cross sections has indicated that the tumors are roughly spherical. These assumptions suggest a model from geometric probability. Consider a medium which contains spheres whose centers are distributed according to a Poisson process in space with constant intensity and whose equatorial radii are distributed according to the cumulative distribution function $F_3(r)$. It is assumed that the tumor number density is small enough so that distinct spheres do not interfere with one another. Now suppose this medium is sliced in a manner independent of the spheres' sizes and locations. Let $F_2(x)$ denote the cumulative distribution function of the (2 dimensional) cross sectional radii from randomly selected slices. The relationship between $F_2$ and $F_3$ was derived by Wicksell (1925), and is

$$F_2(x) = 1 - \frac{1}{\mu} \int_x^R \sqrt{r^2-x^2}\, dF_3(r) \qquad R \geq x \geq 0 \qquad (1.1)$$

where R is an upper bound for the maximum possible value of r and $\mu$ is the mean (3 dimensional) radius,

$$\mu = \int_0^R r\, dF_3(r). \qquad (1.2)$$

Equation (1.1) is obtained by a conditioning argument. If a single sphere
of radius r is cut then the distance from the cutting plane to the center
of the sphere is equally likely to be anywhere between 0 and r and the c.d.f.
of the cross sectional radius is $F_2(x) = 1 - \frac{\sqrt{r^2-x^2}}{r}, 0 \leq x \leq r$. The probability
that a sphere of radius r will be cut is proportional to its radius
times its relative frequency in the sphere population.

In this work we will usually be acting as though we are sampling
from some population of tumors which possess a density $f_3$. The
problem is: Given a sample from $F_2$, obtain a good estimate for the
density $f_3(r) = F_3'(r)$. In practice tumor cross sections can only be
observed if they are larger than some radius $\varepsilon$. In this case, clearly
the experiment does not provide information concerning $f_3(x)$ for $x \leq \varepsilon$.
However, an integral relationship between the two dimensional distribution,
conditional on $x \geq \varepsilon$, and $f_3(x)$ for $x \geq \varepsilon$, can still be obtained. This
was observed by Chover and King (1982) and we give their derivation
below. Let $F_2^{\varepsilon}$ be the conditional distribution of x given $x \geq \varepsilon$.
Defining $\mu_\varepsilon$ by

$$\mu_\varepsilon = \int_\varepsilon^R \sqrt{r^2-\varepsilon^2}\, f_3(r)dr \tag{1.3}$$

it follows from (1.1) that

$$1 - F_2(\varepsilon) = \frac{\mu_\varepsilon}{\mu}, \tag{1.4}$$

hence,

$$1 - F_2^{\varepsilon}(x) = 1 - \frac{F_2(x)-F_2(\varepsilon)}{1-F_2(\varepsilon)} = \frac{\mu}{\mu_\varepsilon} - \frac{\mu}{\mu_\varepsilon} F_2(x). \qquad (1.5)$$

Substituting (1.5) into (1.1) gives

$$F_2^{\varepsilon}(x) = 1 - \frac{1}{\mu_\varepsilon} \int_x^R \sqrt{r^2-x^2} \; f_3(r)dr. \qquad (1.6)$$

The problem now is to estimate $f_3(r)$, $r \geq \varepsilon$ (or, rather, $f_3^{\varepsilon}(r) = f_3(r)/(1-F_3(r))$, given a sample from $F_2^{\varepsilon}$.

The problem of estimating the distribution of sphere sizes in a medium from the cross sections of a randomly oriented slice given a sample cumulative distribution function from $F_2$ is a classical problem in stereology. For the case $\varepsilon = 0$, several approaches have been proposed, including maximum likelihood, regression and nonparametric methods. See Keiding et al. (1972), Nicholson and Merck (1969), Nicholson (1970,1976), Tallis (1970). Recently, Kuk (1982) has placed this problem in the context of estimating a mixing distribution. Watson (1971) discussed the estimation of moments of $f_3$. Anderssen and Jakeman (1975) obtained an estimate of $f_3$ from the inversion formula

$$f_3(r) = \frac{d}{dr} \; \frac{2}{\pi} \int_r^R \frac{dF_2(x)}{\sqrt{x^2-r^2}}$$

They use spectral differentiation and product integration to evaluate the integral. Mendelsohn and Rice (1982) have recently studied a similar problem in which the desired density g and the density h from which observations are made are related by

$$h(r) = \int w(r,x)g(x)dx \qquad (1.7)$$

for a normal kernel w. Their work is somewhat related to the work described here and will be discussed later.

The problem of recovering estimates of $f_3$ from observations on $f_2$ is harder than might appear at first glance because it is ill posed. Here this means that large changes in the true $f_3$ lead to changes in the sample histogram which are imperceptible compared to the sampling error. In particular, "high frequency" components in $f_3$ will not in general be recoverable from medium or even large samples from $F_2^\varepsilon$. For this reason parametric methods (if a parametric form is known) or nonparametric methods which estimate a smooth solution are most likely to be successful. If the true solution is "smooth" then a good nonparametric smoothing method is a promising candidate for recovering the "truth". If the truth is not smooth then such a method should recover the smooth part of $f_3$. Similar remarks have also appeared in Anderssen and Jakeman (1975), Mendelsohn and Rice (1983) and elsewhere, but are worth repeating.

In Section 2.1 we show how the problem of estimating $f_3^\varepsilon$ from a sample from $f_2$ can be converted to the problem of solving an integral equation given noisy data. We can then apply cross validated spline methods for solving ill posed integral equations. These methods have been shown to be successful in a variety of applications. (See Crump and Seinfeld (1982), Merz (1980), Wahba (1977, 1979, 1980, 1982a,b).)

In Section 2.2 we develop a numerical algorithm using certain carefully matched quadrature approximations, which are particularly suited to the application of cross validated spline methods to integral equations like (1.6).

In Section 3 we apply the estimation procedure to a sample of cross sectional mouse liver data obtained by two of us (S.G. and T.P.). The

mouse liver from which this data was taken was exhaustively dissected
and the true distribution of the three dimensional tumors from the subject
mouse determined.  Thus we have a unique opportunity to compare the
estimated distribution with an actual distribution in circumstances
which accurately reflect laboratory experiments.

The results appear to be quite successful.

Convergence properties of this estimate can be obtained by adapting
known techniques for regularized solutions to ill posed linear operator
equations, see, e.g. Cox (1983), Lukas (1981), Silverman (1983), Wahba (1977).
The results will appear in Nychka (1983).  More to the immediate point,
the experimenter would like to know how well the method will recover size
distributions with a sample size and slicing design similar to those encountered
in practice.  We have designed a Monte Carlo experiment to answer this question
for an experiment similar to the laboratory experiment described in Section 3.
This experiment is in the spirit of the recent landmark paper of Diaconis
and Efron (1983).  Some of the results are given in
Section 4.  In general, the accuracy of the estimate is quite impressive,
considering the modest sample size and ill posedness of the problem.
It is, however, difficult to estimate $f_3(r)$ for r near $\varepsilon$ with sample sizes
like those of Section 3.  This is not surprising considering that $f_3$
is subject to length biased  sampling and that large tumors can give
rise to both large and small cross sections.  Thus information in the
data concerning the behavior of $f_3$ near $\varepsilon$ is scanty.  The method described
here extrapolates from data rich  to data poor regions of r in a linear
manner.  In Section 5 we describe how a priori information concerning the
behavior of $f_3$ near $\varepsilon$ can, if available, be incorporated into the estimate.

In Section 6 we show how certain eigensequence plots can provide important insight into the precise degree of ill posedness of this problem, and we discuss the effects of "binning" the data.

In Section 7 some related methods are described, and we describe the very important interplay between statistical smoothing methods, and approximation theoretic methods such as quadrature and finite element methods.

## 2. CROSS VALIDATED SPLINE METHODS FOR ILL POSED LINEAR OPERATOR EQUATIONS

### 2.1 The cross validated spline estimate $f_{\hat{\lambda}}$ for $f_3$

Let $H$ be the (Sobolev) Hilbert space of real-valued functions on $[\varepsilon,R]$,

$$H = \{h: \ h,h' \text{ abs. cont.}, \ h'' \varepsilon L_2[\varepsilon,R]\}.$$

The (usual) model behind cross validated spline methods for integral equations is:

$$z_i = L_i h + \varepsilon_i, \quad i = 1,2,\ldots,n, \tag{2.1}$$

where the $\{\varepsilon_i\}$ are independent zero mean random variables with common unknown variance, and $L_1,\ldots,L_n$ are bounded linear functionals on $H$. See Wahba (1977, 1978, 1980, 1982,a,b). Given data $z = (z_1,\ldots,z_n)^t$, the cross validated spline estimate $h_{\hat{\lambda}}$, for $h$ is obtained as the minimizer in $H$ of

$$\frac{1}{n}\sum_{i=1}^{n}(L_i h - z_i)^2 + \lambda \int_{\varepsilon}^{R}(h''(r))^2 dr \tag{2.2}$$

where the smoothing (bandwidth) parameter is taken as the generalized cross validation (GCV) estimate of $\lambda$. (See Craven and Wahba (1979).)

In the problem under study, let $\hat{F}_2^{\varepsilon}$ be the sample c.d.f. of the cross sectional radii, let $\{P_i\}_{i=1}^{n}$ be a partition of the interval $[\varepsilon,R]$, $\varepsilon = P_1 < P_2 < \ldots < P_n < R$, and let $z_i$ be the fraction of all observations in the ith bin, $[P_i, P_{i+1})$.

Then

$$z_i = \hat{F}_2^{\varepsilon}(P_{i+1}) - \hat{F}_2^{\varepsilon}(P_i) = F_2^{\varepsilon}(P_{i+1}) - F_2^{\varepsilon}(P_i) + \varepsilon_i, \tag{2.3}$$

where the $\varepsilon_i$ are random variables. If the observations are an independent

sample from $F_2^\varepsilon$, then the $\{\varepsilon_i\}$ will have zero mean and be jointly asymptotically normal and only weakly correlated. In this work we are going to ignore the fact that the variances of the $\varepsilon_i$ are not necessarily the same. (Various reweighting schemes are possible, see Cox (1970), Villalobos and Wahba (1983).) Letting $h = f_3/\mu_\varepsilon$, and setting

$$L_i h = \int_{P_i}^{R} \overline{\sqrt{r^2 - P_i^2}} h(r) dr - \int_{P_{i+1}}^{R} \overline{\sqrt{r^2 - P_{i+1}^2}} h(r) dr \qquad (2.4)$$

(2.3) becomes (with the aid of (1.6)),

$$z_i = L_i h + \varepsilon_i.$$

Given z, we let $h_\lambda$ be the minimizer of (2.2) in $H$, and let $f_\lambda$ be

$$f_\lambda(r) = h_\lambda(r) / \int_\varepsilon^R h_\lambda(s) ds. \qquad (2.5)$$

Our estimate $\hat{f}_3$ is then $f_{\hat\lambda}$, where $\hat\lambda$ is the GCV estimate of $\lambda$. (Note that $\int_\varepsilon^R h_{\hat\lambda}(s) ds$ is an estimate for $1/\mu_\varepsilon$ .). The estimate obviously integrates to 1, but it is not required to be positive. Negativity was not a problem with the actual mouse liver data. In one of the Monte Carlo examples the estimate went negative and we have truncated the estimate in the plots. If desired, non negativity constraints can be added to the problem of (2.2), see Wahba (1982a), Villalobos (1983).

## 2.2 The numerical method for computing $f_{\hat{\lambda}}$

Using known but scattered results, we next give an efficient numerical procedure for computing (a very good approximation to) the minimizer of (2.2) and the GCV estimate $\hat{\lambda}$ of $\lambda$. The method is readily implemented for n less than a few hundred. In all our calculations, n will be 80, and the bins are equally spaced in log x between $\varepsilon$ and R. The log spacing is a crude variance stablizing spacing for our mouse liver data. For the actual and most of the Monte Carlo data the number of observed cross sections was between 150 and 450. The choice of n = 80 bins is large enough so that the binning is not doing any appreciable smoothing. Binning as smoothing will be discussed in further detail in Section 6.

As with any ill posed problem, care must be taken in the actual calculation of the solution, or garbage may result from dividing random or roundoff errors by small eigenvalues. It will be seen here and in Sections 6 and 7 that the numerical analysis and the estimation procedure can become inextricably intertwined in ill posed problems. Approximation theoretic methods become smoothing procedures and vice versa. For completeness, and to allow discussion of this point, we outline the major steps of our numerical method here, pointing out the steps developed particularly for the problem at hand.

Using the results in Kimeldorf and Wahba (1971), Wahba (1978) and Wahba and Wendelberger (1980) an explicit formula for $h_\lambda$, the minimizer of (2.2) in $H$ can be given as follows. Under the inner product

$$\langle f,g\rangle_H = f(\varepsilon)g(\varepsilon) + f'(\varepsilon)g'(\varepsilon) + \int_\varepsilon^R f''(r)g''(r)dr$$

$H$ is a reproducing kernel Hilbert space. The reproducing kernel for $H$ with this inner product is

$$Q(r,s) = 1 + (r-\varepsilon)(s-\varepsilon) + Q_1(r,s), \quad \varepsilon \le r, s \le R \qquad (2.5)$$

where

$$Q_1(r,s) = \frac{(r-\varepsilon)^2(s-\varepsilon)}{2} - \frac{(r-\varepsilon)^3}{6} \quad r \le s$$

$$\frac{(r-\varepsilon)(s-\varepsilon)^2}{2} - \frac{(s-\varepsilon)^3}{6} \quad r \ge s$$

Let $\phi_1(r) = 1$, $\phi_2(r) = (r-\varepsilon)$, and $\xi_i(r) = L_i(Q_1(\cdot,r))$, where $L_i$ is given by (2.4) and $L_i(Q_1(\cdot,r))$ means that $L_i$ is applied to $Q_1(s,r)$ considered as a function of s. Let T be the n×2 matrix with iνth entry $\tau_{i\nu} = L_i\phi_\nu$, $\nu = 1,2$, and let K be the n×n matrix with ijth entry $k_{ij} = \int \xi_i''(r)\xi_j''(r)dr$, If T is of rank 2, $h_\lambda$ is uniquely determined, and given by

$$h_\lambda(r) = \sum_{i=1}^n c_i\xi_i(r) + \sum_{\nu=1}^2 d_\nu\phi_\nu(r) \qquad (2.6)$$

where $c = (c_1,\ldots,c_n)'$ and $d = (d_1,d_2)'$ satisfy

$$(K+n\lambda I)c + Td = z \qquad (2.7a)$$

$$T'c = 0 \qquad (2.7b)$$

The GCV estimate $\hat{\lambda}$ of $\lambda$ is the minimizer of the cross validation function $V(\lambda)$,

$$V(\lambda) = 1/n(||(I-A(\lambda))z||^2)/(\frac{1}{n}Tr(I-A(\lambda)))^2, \qquad (2.8)$$

where $A(\lambda)$ is the n×n "influence matrix" defined by

$$\begin{bmatrix} L_1(h_\lambda) \\ \vdots \\ L_n(h_\lambda) \end{bmatrix} = A(\lambda)z$$

From, e.g. Wahba and Wendelberger (1980) it is known that

$$I-A(\lambda) = Q(QKQ'+n\lambda I)^{-1}Q', \qquad (2.9)$$

where $Q$ can be taken as any n×n-2 matrix whose n-2 columns are linearly independent and perpendicular to the 2 columns of $T$. The numerical problem now is to compute the minimizer $\hat{\lambda}$ of $V(\lambda)$, and $h_{\hat{\lambda}}$.

In this problem closed form expressions can be obtained for the $\{\xi_i\}$ and $\{\tau_{i\nu}\}$, and are given in Appendix A. Unfortunately we were unable to find a closed form expression for $k_{ij} = \int_\varepsilon^R \xi_i''(r)\xi_j''(r)dr$, so some form of quadrature must be used. It is not at all clear that just applying the nearest handy quadrature formula to obtain approximations to the entries $k_{ij}$ of K is appropriate. In particular, the nonnegative definiteness of K could easily be lost, leading to problems in the calculation of $\hat{\lambda}$, see below.

The following form of "matched quadrature" can be used to avoid this problem. The particular form of "matched quadrature" chosen is motivated by a) the fact that $\xi_i^{(\nu)}(\varepsilon) = 0$, $\nu = 0,1$, $i = 1,2,\ldots,n$, and b) the desire to do as little quadrature approximation as possible by exploiting the known closed form expressions for $\xi_i$ and $\tau_{i\nu}$.

First, let $\eta_i$ be the representer of $L_i$ in $H$, that is,

$$L_i h = <\eta_i, h>.$$

It is known that $\eta_i = \xi_i + a_{i1}\phi_1 + a_{i2}\phi_2$ for some $a_{i1}, a_{i2}$, see, e.g.
Kimeldorf and Wahba (1971). Note that (2.2) may be rewritten

$$\frac{1}{n}\sum_{i=1}^{n}(<\eta_i,h>-z_i)^2 + \lambda\int_\epsilon^R (h''(r))^2 dr.$$

Now, choose a fine grid of N+1 points, $\epsilon = s_0 < s_1 < s_2 < ... < s_N = R$
and, for any h let $P_N h$ be that element in $H$ which minimizes $J(h)$
subject to $(P_N h)(s_\ell) = h(s_\ell)$, $\ell = 0,1,2,...,N$, and $(P_N h)'(\epsilon) = h'(\epsilon)$.
$P_N h$ will be a cubic interpolating spline subject to the left boundary
conditions. The "matched quadrature" consists of approximating
$L_i$ by $\tilde{L}_i$, where $\tilde{L}_i$ is the linear functional on $H$ defined by

$$\tilde{L}_i h = <P_N \eta_i, h>.$$

We are now in a position to solve the approximate problem:

Minimize

$$\frac{1}{n}\sum_{i=1}^{n}(\tilde{L}_i h - z_i)^2 + \lambda\int_\epsilon^R (h''(r))^2 dr, \tag{2.11}$$

in $H$, <u>exactly</u>. This is easily done using the formula (2.6) since, it
can be shown that $\tilde{\tau}_{i\nu} \equiv \tilde{L}_i \phi_\nu = L_i \phi_\nu = \tau_{i\nu}$ and $\tilde{L}_i(Q_1(\cdot,r)) = P_N\xi_i = \tilde{\xi}_i$, say. [1]
$\tilde{k}_{ij} = \int_\epsilon^R \tilde{\xi}_i''(r)\tilde{\xi}_j''(r)dr$ is readily evaluated exactly, since the $\{\tilde{\xi}_i\}$ are
piecewise polynomials. The minimizer $h_\lambda$ of (2.11) is given by (2.6)
and (2.7) with $\xi_i$ replaced by $\tilde{\xi}_i$ and K replaced by $\tilde{K} = \{\tilde{k}_{ij}\}$. The cross
validation function $\tilde{V}(\lambda)$ for this problem is given by (2.8) and (2.9) with $A(\lambda)$

---

[1] The procedure we used for computing $P_N\xi_i$ is given in Appendix B.

replaced by $\tilde{A}(\lambda)$ defined by replacing K by $\tilde{K}$ in (2.9). $Q\tilde{K}Q'$ will be nonnegative definite.

Given $\tilde{K}$ and T, we give an efficient procedure for minimizing $\tilde{V}(\lambda)$ and computing c and d.

1. Use LINPACK (Dongarra, et. al. (1979)) to find the QR decomposition of T, to obtain

$$T = (Q_1 : Q_2) \begin{pmatrix} R_1 \\ \cdots \\ 0 \end{pmatrix}$$

where $Q_2$ is an n×n-2 matrix with $Q_2'Q_2 = I_{n-2 \times n-2}$, $Q_2'T = 0$ and $R_1$ is upper triangular. The Q appearing in (2.9) can be taken as $Q_2$.

2. Let $B = Q_2'\tilde{K}Q_2$ and use EISPACK (Smith et.al.(1976) to find the eigenvalue eigenvector decompositon $UD_BU'$ of B, where $b_\nu^2$, $\nu = 1,2,\ldots,n-2$ are the n-2 diagonal entries of $D_B$ (eigenvalues of B), and the n-2 columns of U are the eigenvectors of B. Then

$$\text{Trace}(I-\tilde{A}(\lambda)) = \sum_{\nu=1}^{n-2} \frac{n\lambda}{b_\nu^2+n\lambda}$$

$$(I-\tilde{A}(\lambda))z = n\lambda Q_2 U(D_B+n\lambda)^{-1}U'Q_2'z.$$

3. Letting $w = U'Q_2'z$, then

$$\tilde{V}(\lambda) = \frac{1}{n}\sum_{\nu=1}^{n-2}\left(\frac{n\lambda w_\nu}{b_\nu^2+n\lambda}\right)^2 / \left(\frac{1}{n}\sum_{\nu=1}^{n-2}\frac{n\lambda}{b_\nu^2+n\lambda}\right)^2$$

$$c = Q_2 U(D_B+n\lambda I)^{-1}w$$

and d is obtained by solving

$$R_1 d = Q_1'(z-\tilde{K}c).$$

$\tilde{V}(\lambda)$ is minimized by a global search in log $\lambda$. If $n\lambda$ is much smaller than the smallest $b_\nu{}^2$ or much larger than the largest $b_\nu^2$, it may be taken as 0 or $\infty$, respectively, so this limits the region required to be searched. $\int_\epsilon^R h_{\hat{\lambda}}(r)dr$ is easily evaluated.

It is useful to note that if $h_{\hat{\lambda}}$ is to be obtained for repeated samples, with the same bins, the cost is quite modest for runs after the first, since the expensive calculations involve the calculation of $Q_1, Q_2, R$, U and $D_B$ and these need only be computed once as they do not depend on the data.

The above procedure appears in Wendelberger (1981), and has been found to work well in similar problems for n as large as 350. In the calculations that follow we used N = 80, with the $s_i$ equally spaced. Further discussion of the choice of N appears in Section 6.

## 3. NUMERICAL RESULTS WITH THE LABORATORY DATA

The liver being sliced fits roughly into a box about 7500 × 7500 microns (μ) square by 2380μ deep (100μ=.01cm.), and, for the experimental data studied is sliced perpendicular to the short dimension in 21 equally spaced slices (of negligible thickness) 50 microns apart, through the central 1000 microns of the block, to be called the slicing region. Figure 3.1 gives a schematic diagram of the slicing design.
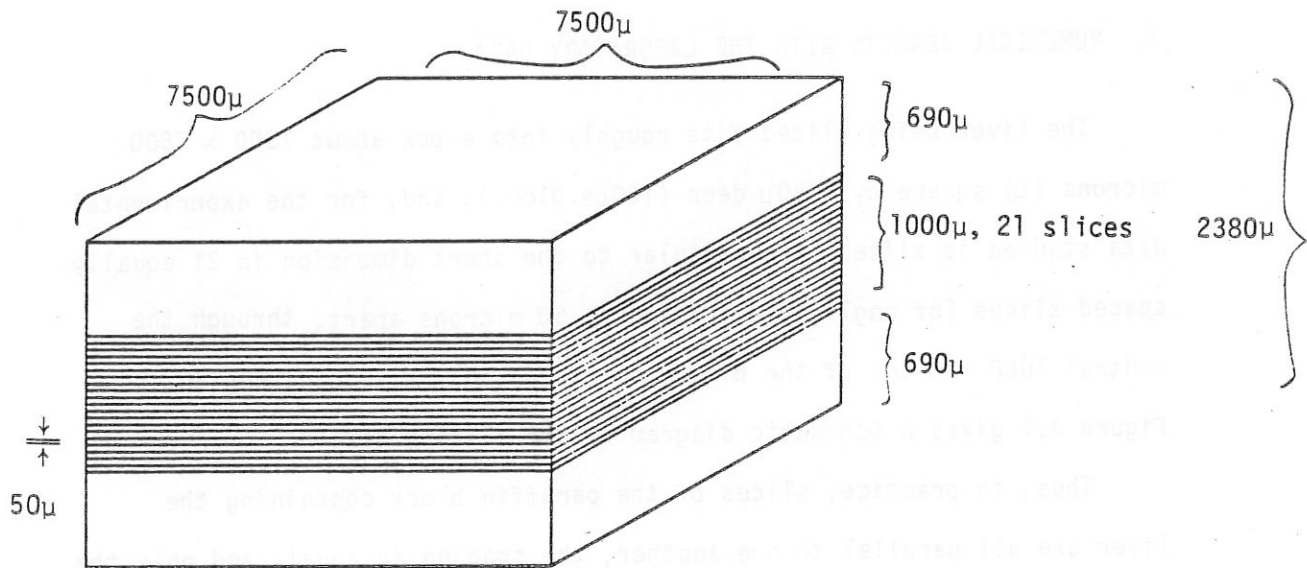
Thus, in practice, slices of the paraffin block containing the liver are all parallel to one another, the spacing is equal, and only the "phase" of the tumors with respect to the slicing grid is random. In this experiment, it was decided in the laboratory that $\varepsilon$ = 38.36 [1] microns was the smallest cross sectional radius reliably detected by all the personnel identifying cross sections. Smaller cross sections, when observed were ignored. Spherical tumors of three dimensional radius greater than 45.78 microns and lying wholly in the slicing region will be observed in at least one slice, while tumors with radius between $\varepsilon$ and 45.78 may or may not be observed. Figure 3.2 gives the probability that a sphere of radius r which lies wholly in the slicing region will have at least one observed cross section. Equation (1.1) still holds, but a little reflection will show that if the spacing is uniform, and spheres can be sliced more than once, the sampling variance will become smaller as the spacing becomes finer.

With this slicing 154 tumor cross sections were observed. Figure 3.3 gives a histogram of the observed cross sectional radii using the bins $[P_i, P_{i+1})$. Figure 3.4
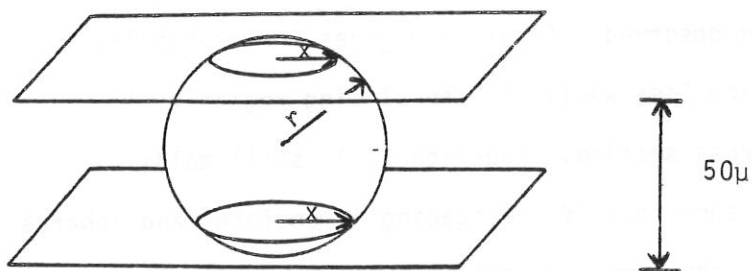
---

[1] The observations do not actually have 4 figure accuracy. we are ignoring this.

(a) slicing design



(b.) Detail of sphere intersected by two slices
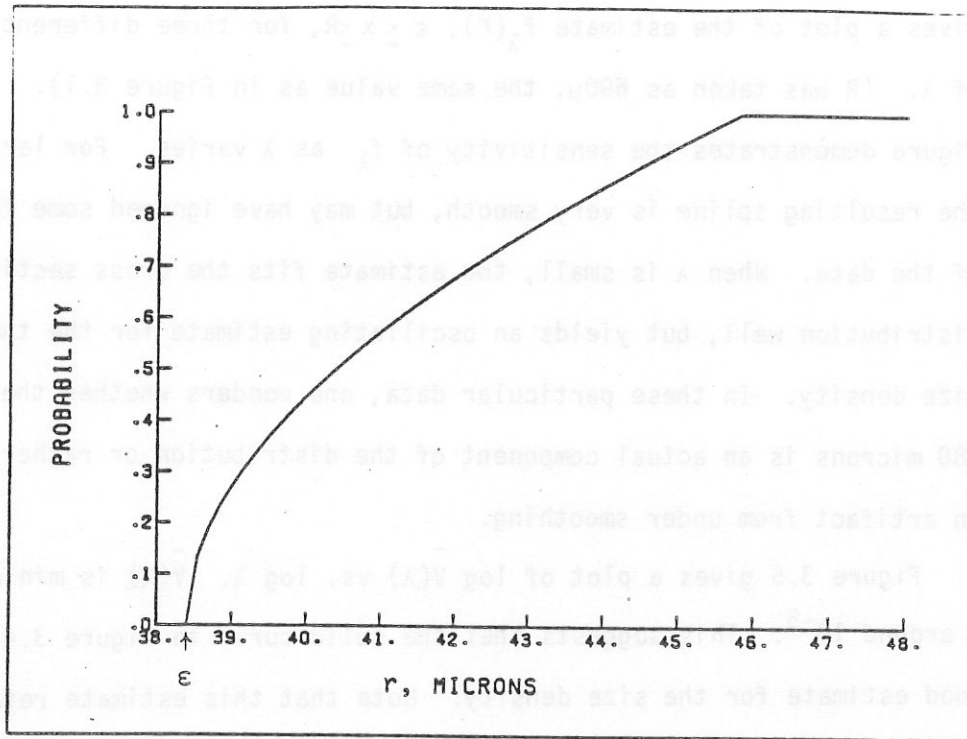
Figure 3.1  Schematic diagram of the slicing design

Figure 3.2   Probability that a sphere of radius r
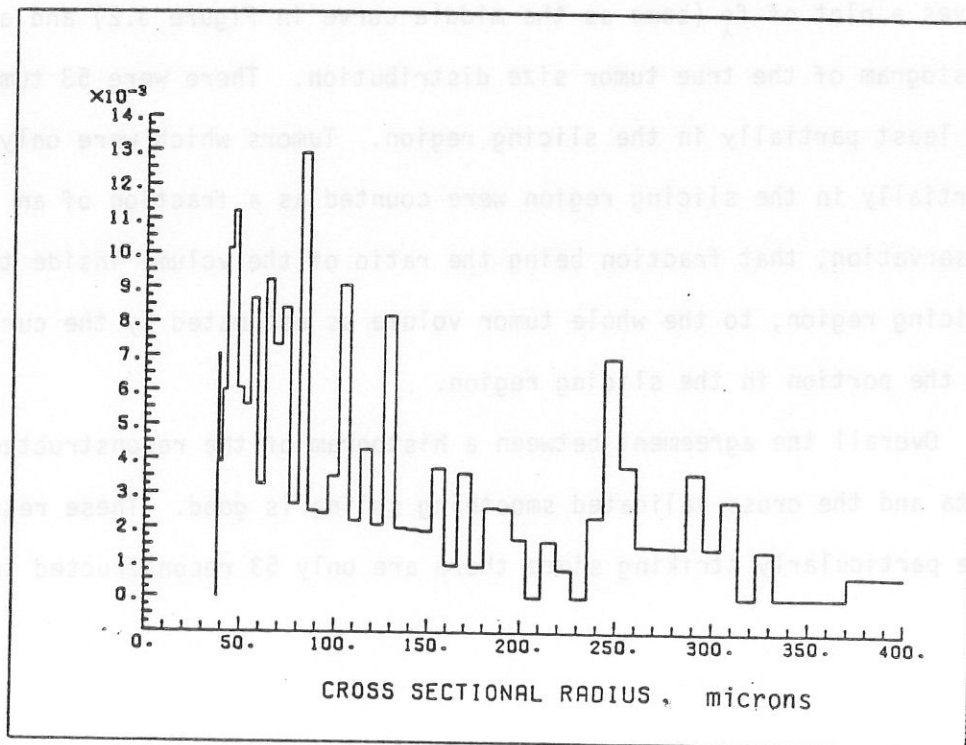will have at least one observed cross section.



Figure 3.3   Histogram of cross sectional radii,
154 observed cross sections.

gives a plot of the estimate $f_\lambda(r)$, $\varepsilon \le x \le R$, for three different values
of $\lambda$. (R was taken as 690μ, the same value as in Figure 3.1). This
figure demonstrates the sensitivity of $f_\lambda$ as $\lambda$ varies. For large $\lambda$,
the resulting spline is very smooth, but may have ignored some features
of the data. When $\lambda$ is small, the estimate fits the cross sectional
distribution well, but yields an oscillating estimate for the tumor
size density. In these particular data, one wonders whether the mode at
280 microns is an actual component of the distribution or rather just
an artifact from under smoothing.

Figure 3.5 gives a plot of log $\tilde{V}(\lambda)$ vs. log $\lambda$. $\tilde{V}(\lambda)$ is minimized for
$\lambda$ around $10^{-5}$. This suggests that the solid curve in Figure 3.4 is a
good estimate for the size density. Note that this estimate retains
a mode at around 280 microns. To compare $f_{\hat{\lambda}}$ with the true $f_3^\varepsilon$, the slicing
region was completely dissected by very fine slicing. Figure 3.6
gives a plot of $f_{\hat{\lambda}}$ (same as the middle curve in Figure 3.2) and a
histogram of the true tumor size distribution. There were 53 tumors
at least partially in the slicing region. Tumors which were only
partially in the slicing region were counted as a fraction of an
observation, that fraction being the ratio of the volume inside the
slicing region, to the whole tumor volume as estimated by the curvature
of the portion in the slicing region.

Overall the agreement between a histogram of the reconstructed
data and the cross validated smoothing spline is good. These results
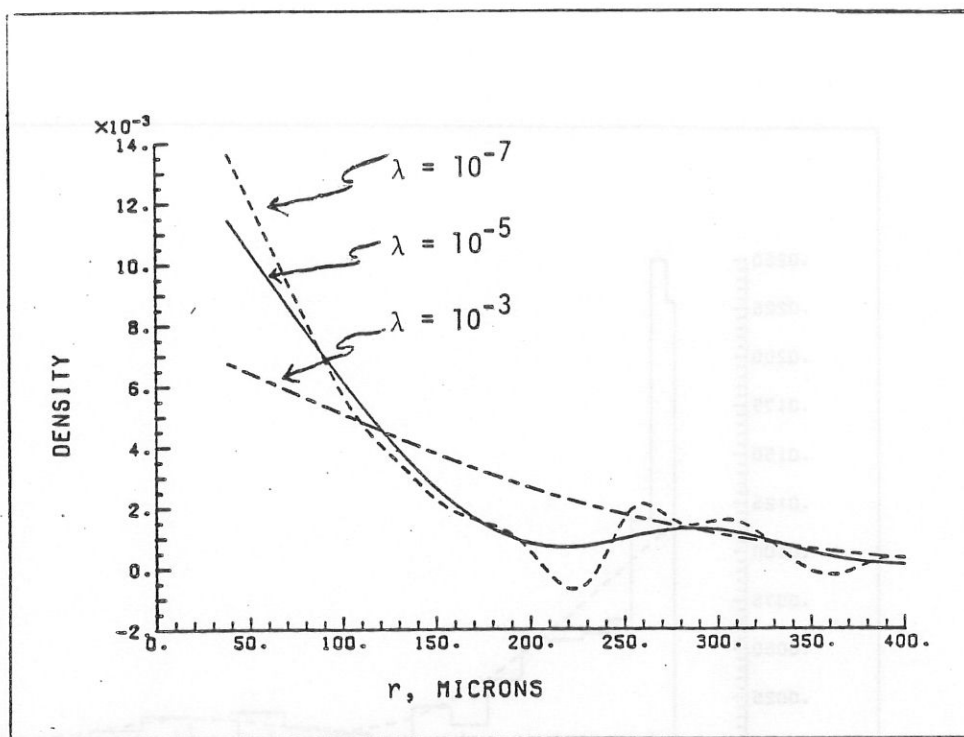are particularly striking since there are only 53 reconstructed tumors

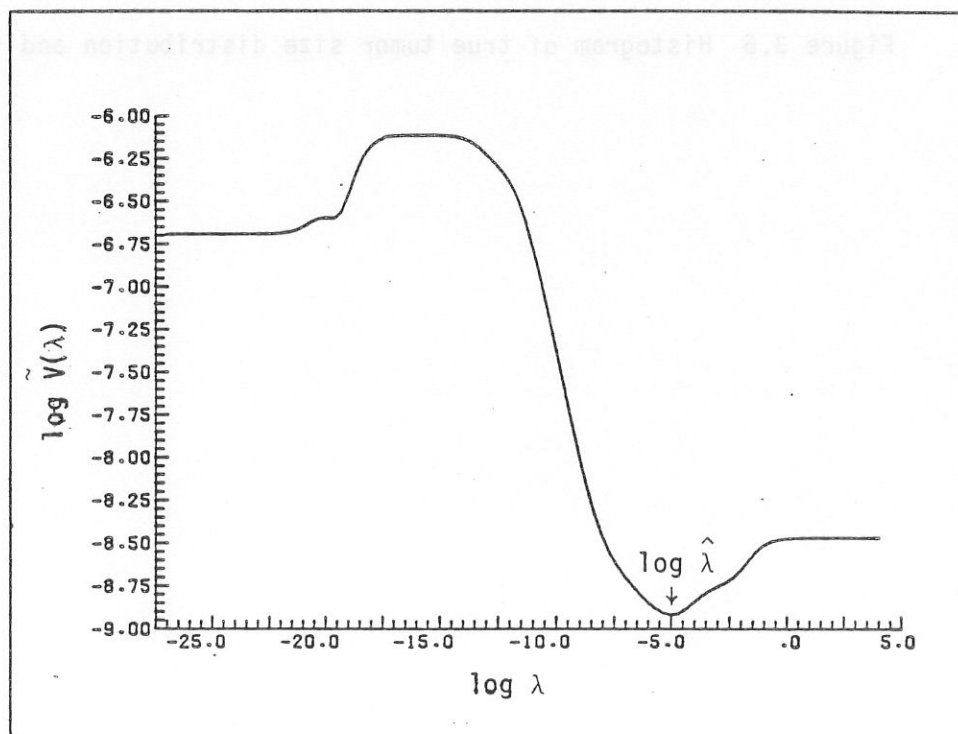Figure 3.4   The estimate $f_\lambda(r)$, for 3 different values of $\lambda$.



Figure 3.5   The cross validation function $\tilde{V}(\lambda)$.
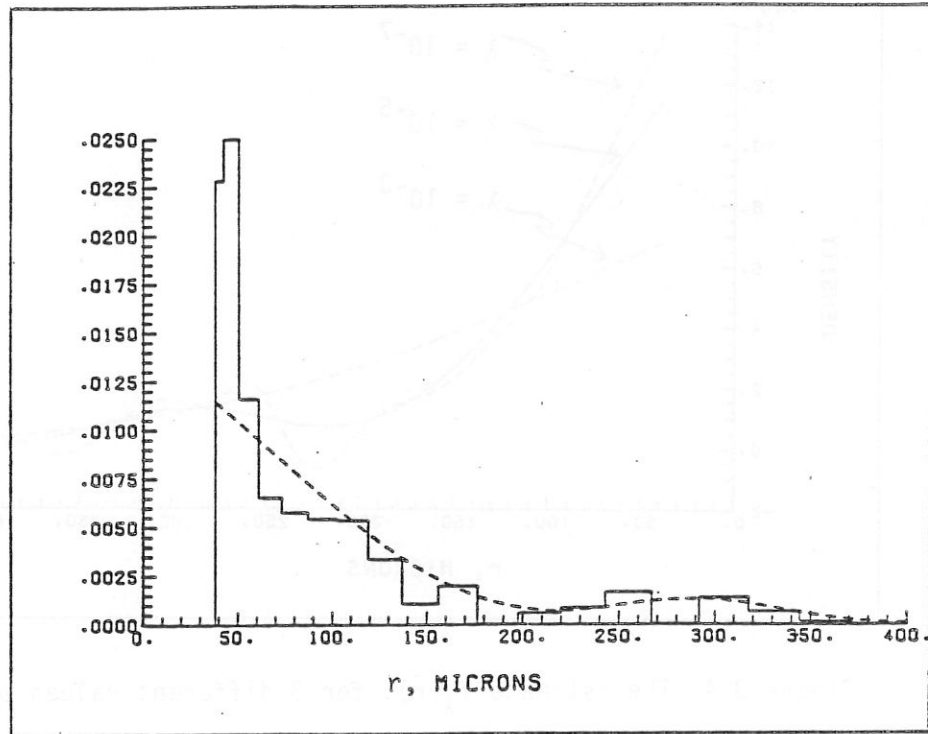
Figure 3.6  Histogram of true tumor size distribution and $f_{\hat{\lambda}}$.

in the tissue sample, although, of course, the systematic sampling helps.
The concentration of tumors around 280 microns predicted by the spline
is an actual feature of the reconstructed data.  However, close to the
lower limit, ε, the spline underestimates the reconstructed distribution.

## 4. MONTE CARLO EXPERIMENTS

We studied the sampling properties of the estimate by Monte Carlo methods designed to mimic the effects of multiple sampling of large tumors, as well as edge effects, as they actually occur in the mouse experiment.

The geometry of the Monte Carlo experiment is exactly that described in Figure 3.1, where, however, R in that figure may take on other values. A pseudo random number for the total number of spheres was generated according to a Poisson distribution with mean equal to the volume of the entire block × 900 tumors/cc. (The actual mouse had a tumor number density of about 900/cc.) If the number of spheres is $n_3$, then $n_3$ "centers" are uniformly distributed throughout the entire block. For each center, a random radius was generated according to the density $f_3$. Twenty one parallel, infinitely thin slices 50 microns apart were then made through the shaded region and the radii of all (two dimensional) intersections greater than $\varepsilon$ were recorded.

There are now at least two ways of defining the "true" distribution of the three dimensional radii in this experiment. One is as the "theoretical" distribution determined by the density $f_3$ from which the pseudo random radii were drawn. The second is as the "actual" distribution of the three dimensional radii that were actually drawn. For comparison purposes we will display both the "theoretical" density and a histogram of the "actual" distribution as defined above. (The "actual" distribution is defined here a little differently than the "true" distribution of Section 3, since tumors in the block but outside the slicing region can be counted.)

Experimenters will likely want to focus on the "actual" distribution if they are interested in a single mouse, and on the "theoretical" distribution if they consider a single mouse as a member of some "super-population".

We present the results of four Monte Carlo studies. In each of the studies six replicates were performed. A replicate consists of drawing a sample of tumors, slicing the block, recording the observed cross sectional radii and computing the estimate $\hat{f}_\lambda$.

Experiment 1 was very roughly designed to mimic the number density and theoretical size density $f_3$ of six experimental mice, one of which has been described in Section 3. The theoretical $f_3$ was taken as a Weibull density that approximates the data of Figure 3.6. R in this experiment was 690 microns. (R in the definition of $J(\cdot)$ and in Figure 3.1 have been taken to be the same.) The number of tumors/replicate in the entire block averaged 115 with about 49% of them having recorded intersections. The number of observed cross sections averaged 204. Figure 4.1 shows the results of the six replications. The solid curve in the upper left plot is the theoretical Weibull curve, the histograms represent the "actual" size distributions, and the dashed lines are the estimates $\hat{f}_\lambda$. While the overall shape of the estimate is quite good, a tendency to underestimate the density near the cutoff is evident in 4 of the six replicates. Experiment 2 studies a density with different behavior near $\varepsilon$. $f_3$ is a truncated Beta density. The average number of tumors in the block was 113 of which around 60% had recorded intersections, the average number of observed cross sectional radii was 426. In this experiment most replicates
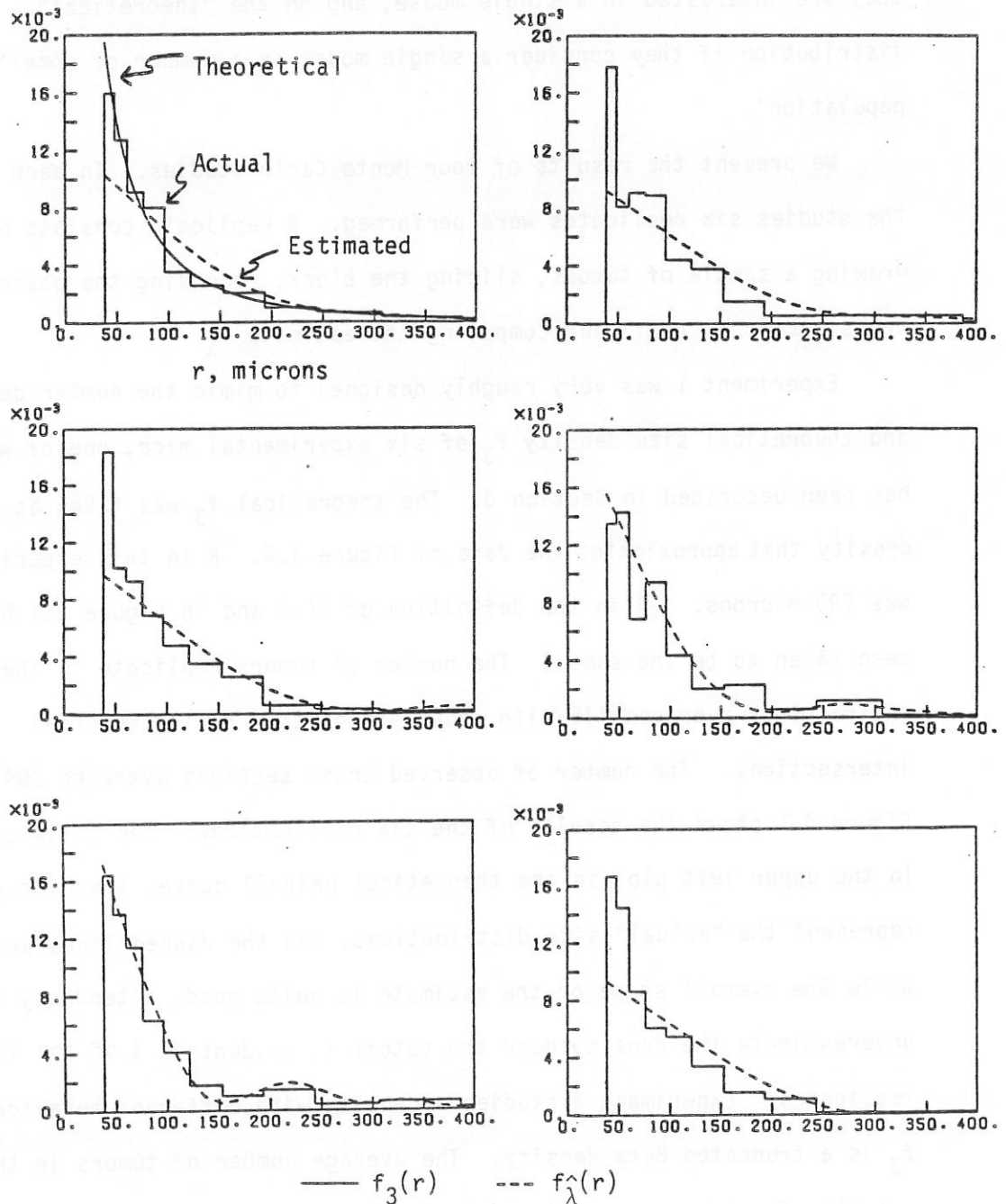
Figure 4.1    Experiment 1.   Theoretical density, histogram
for actual distribution and $f_{\hat{\lambda}}$, six replicates.

Weibull theoretical density.

r, microns

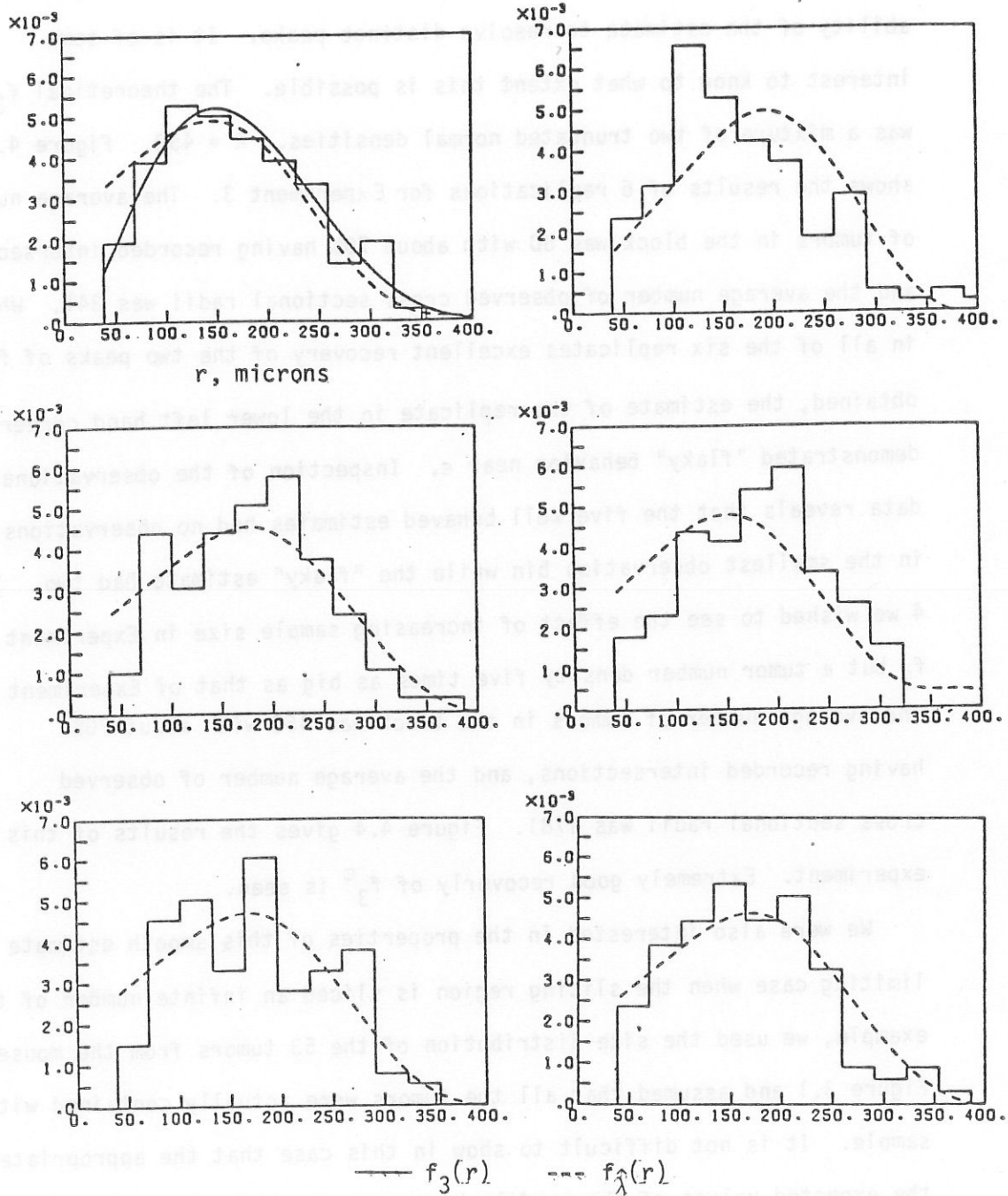$$\underline{\hspace{1cm}} \ f_3(r) \qquad \text{-- --} \ f_{\hat{\lambda}}(r)$$

Figure 4.2   Experiment 2.  Theoretical density,
observed distribution and $f_{\hat{\lambda}}$, six replicates.
Truncated Beta theoretical density,

overestimated $f_3$ near $\varepsilon$ while, overall, the shape of $f_3$ is quite good,
particularly for larger r. In Experiment 3 we wished to examine the
ability of the estimate to resolve distinct peaks. It is of some
interest to know to what extent this is possible. The theoretical $f_3$
was a mixture of two truncated normal densities. R = 450. Figure 4.3
shows the results of 6 replications for Experiment 3. The average number
of tumors in the block was 80 with about 70% having recorded intersections
and the average number of observed cross sectional radii was 341. While
in all of the six replicates excellent recovery of the two peaks of $f_3$ was
obtained, the estimate of the replicate in the lower left hand corner
demonstrated "flaky" behavior near $\varepsilon$. Inspection of the observational
data reveals that the five well behaved estimates had no observations
in the smallest observation bin while the "flaky" estimate had two. In Experiment
4 we wished to see the effect of increasing sample size in Experiment 3. The same
$f_3$ but a tumor number density five times as big as that of Experiment 3 was used.
The average number of tumors in the block was 460 with about 70%
having recorded intersections, and the average number of observed
cross sectional radii was 1781. Figure 4.4 gives the results of this
experiment. Extremely good recoverly of $f_3^\varepsilon$ is seen.

We were also interested in the properties of this smooth estimate in the
limiting case when the slicing region is sliced an infinte number of times. For an
example, we used the size distribution of the 53 tumors from the mouse of
Figure 3.1 and assumed that all the tumors were actually contained within the
sample. It is not difficult to show in this case that the appropriate data is
the expected values of the profile histogram for this discrete distribution (Nychka
(1983). A plot of $f_{\hat{\lambda}}$ appears in Figure 4.5, superimposed on a histogram
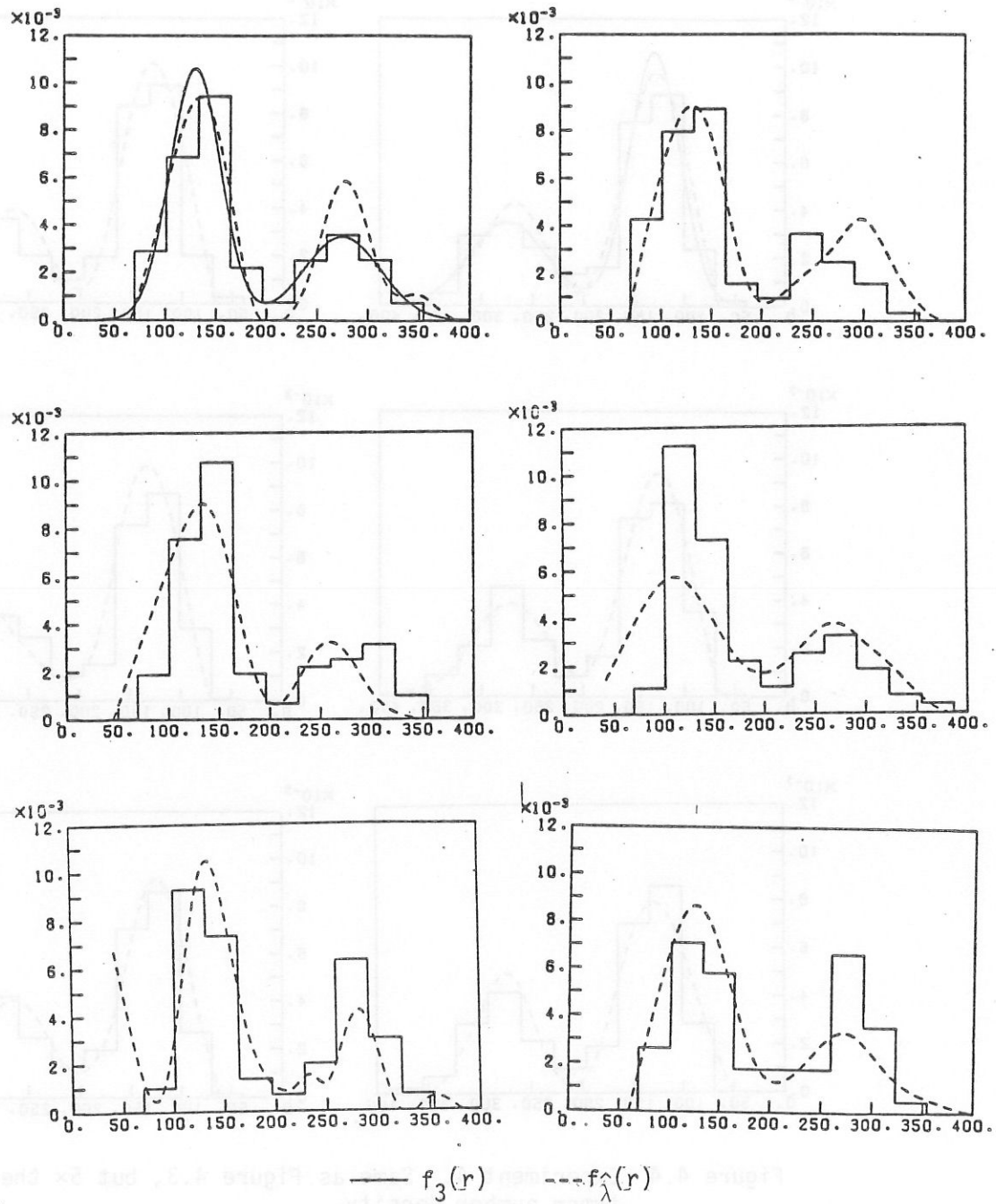of the "theoretical" distribution.

Figure 4.3  Experiment 3.  Theoretical density, histogram for actual
distribution, and $f_{\hat{\lambda}}$, six replicates.  Truncated normal
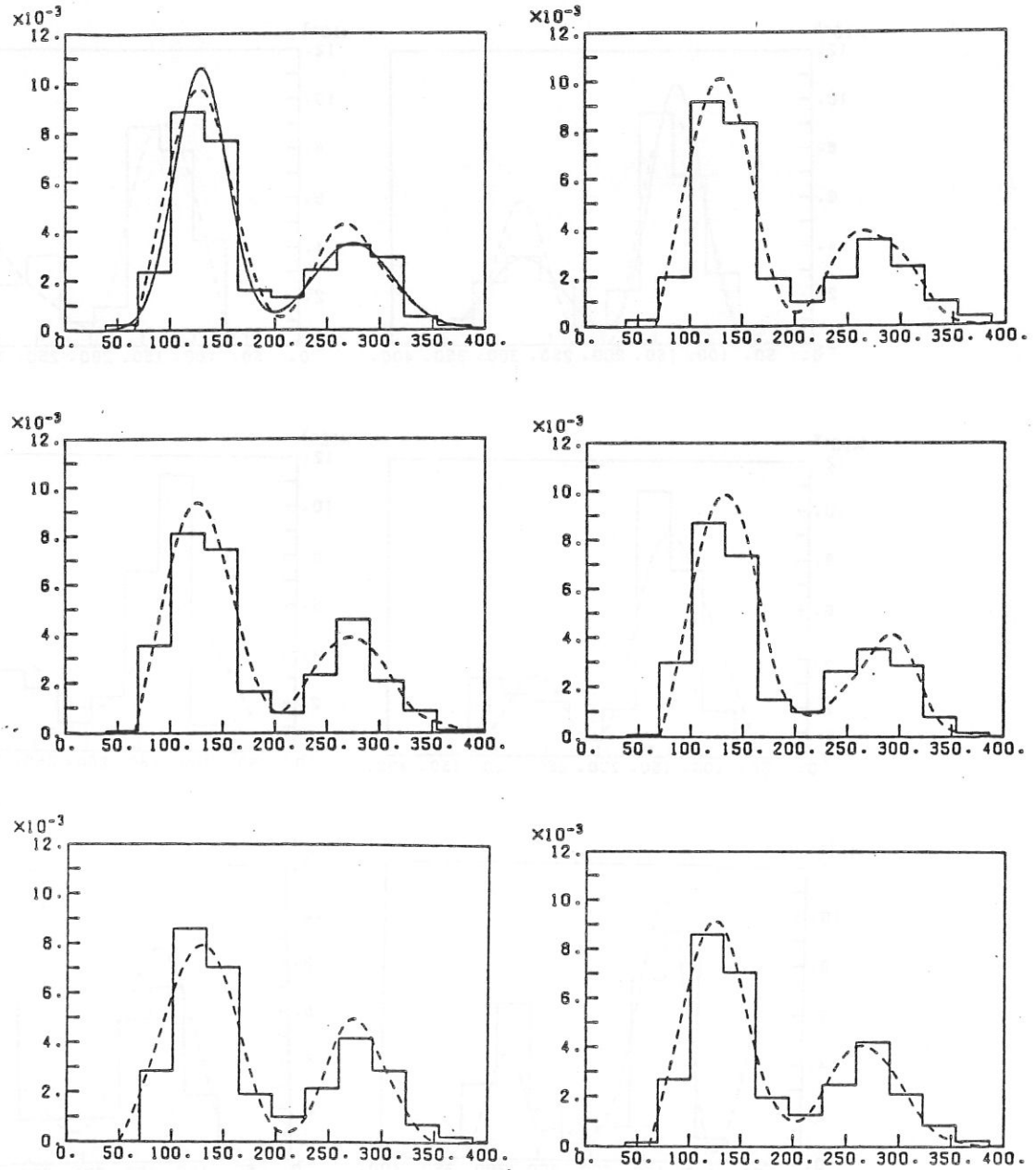mixture theoretical density.

Figure 4.4   Experiment 4.   Same as Figure 4.3, but 5× theoretical
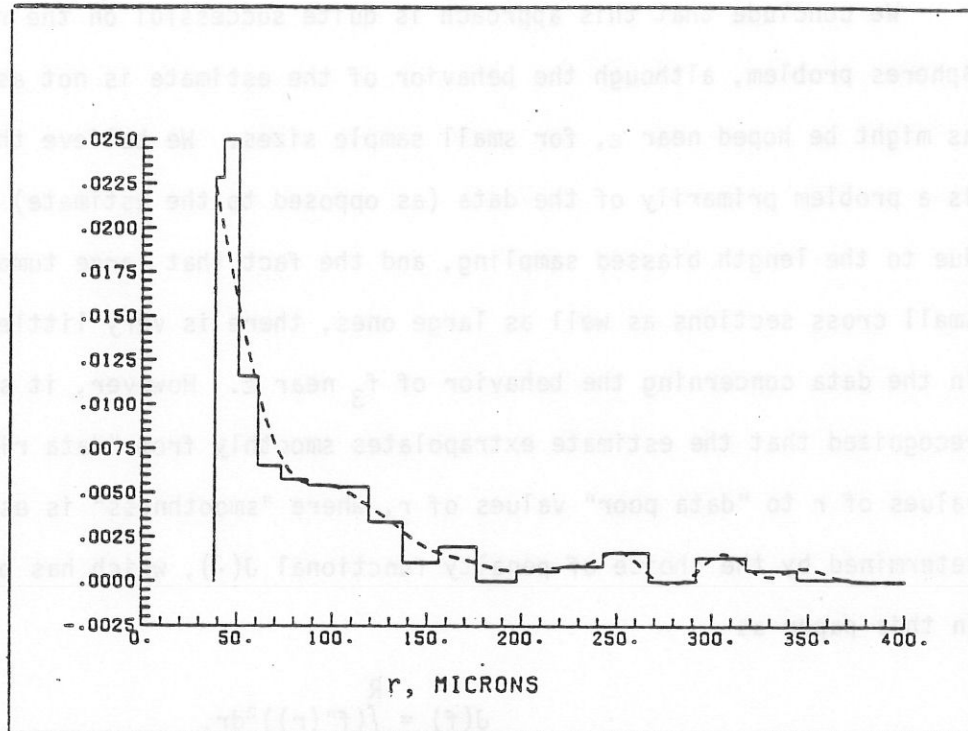tumor number density.

Figure 4.5  Smooth estimate with an infinite sample size
from a discrete theoretical distribution.

## 5. MODIFICATIONS FOR END EFFECTS

We conclude that this approach is quite successful on the random spheres problem, although the behavior of the estimate is not as good as might be hoped near $\varepsilon$, for small sample sizes. We believe that this is a problem primarily of the data (as opposed to the estimate) since, due to the length biassed sampling, and the fact that large tumors have small cross sections as well as large ones, there is very little information in the data concerning the behavior of $f_3$ near $\varepsilon$. However, it should be recognized that the estimate extrapolates smoothly from "data rich" values of r to "data poor" values of r, where "smoothness" is essentially determined by the choice of penalty functional $J(\cdot)$, which has been taken in this paper as

$$J(f) = \int_{\varepsilon}^{R} (f''(r))^2 dr.$$

The null space of $J(\cdot)$ is the linear functions, thus, where there is insufficient information in the data extrapolation will be linear. In this problem, the penalty functional could have been replaced by, say,

$$J(f) = \int_{\varepsilon}^{R} (f'''(r))^2 dr$$

in which case, the extrapolation would have been quadratic. If prior information concerning the behavior of $f_3$ near $\varepsilon$ were available from some external source, then this information could be included in the cross validated spline estimate by appropriately modifying $J(\cdot)$. For

example, suppose it was known that $f_3$ behaves like a particular
negative exponential density g, say.  This information may be incorporated
in the estimate by replacing

$$J(f) = \int^R (f''(r))^2 dr$$

by, for example

$$J(f) = ||P_g f||^2_Q$$

where $P_g f$ is the projection of $f$ onto the orthcomplement of span
$\{g, \phi_1, \phi_2\}$.  Then extrapolation from data rich to data poor regions (i.e. near
$\varepsilon$) will proceed via Bayesian information that the true f has negative
exponential behavior there.[2]  The abstract idea behind this approach
may be found in Wahba (1978), Section 3.  For details of the application
to this problem, see Nychka (1983).  See also the remarks in Silverman
(1982), where a normal density is in the null space of his penalty functional.
For a different approach to modifying boundary behavior, see Gasser
and Muller (1979).

---

[2] We are compelled to report, however, that for the mouse liver problem
behavior of $f_3^\varepsilon$ near $\varepsilon$ was something of a surprise.

6. QUANTITATIVE ILL POSEDNESS, EIGENSEQUENCE PLOTS, AND THE CHOICE OF n

Since the cost of the numerical calculations increase rapidly with n (for the first estimate computed), it is tempting to choose n fairly small. If n is much less than the number of observations, it may act as a smoothing parameter. Using n as a smoothing parameter can be justified theoretically, from an asymptotic point of view (for example, see Wahba (1975)). However, it is our numerical experience that when there is a relatively small amount of information about the solution available in the data, then smoothing by binning can result in loss of fine structure in the estimate that would be observable if $\lambda$ were allowed to do most of the smoothing. Thus, we set out in this problem to choose n large enough so that little or no smoothing is being done at the binning step.

However, since this problem is ill posed, increasing n beyond some point will not retain much more information, even if the sample size were infinite.

Inspection of the computed eigenvalues $b_\nu^2$, $\nu = 1,2,\ldots,n-2$ can be a valuable procedure in studying this question and we describe how below. First, given the bins, let $K_n$ be the operator with domain $H$ and range $E_n$ which maps f to $K_n f = (L_1 f,\ldots,L_n f)$. Then $K_n$ is analogous to the design matrix X in the usual regression problem $y = X\beta + \epsilon$, and the role of XX' is played by the n×n gram matrix $\Sigma$ with ijth entry $\langle \eta_i, \eta_j \rangle$. Inspection of the eigenvalues of $\Sigma$ thus provides important information on the effective dimension of the range of $K_n$, when the domain of $K_n$ is $H$.

In some ill posed problems, $\Sigma$ is theoretically of full rank but has fewer than n eigenvalues which are actually larger than machine double precision 0. For an extreme example, see Wahba (1979). Now, since $\eta_i = \xi_i + a_{i1}\phi_1 + a_{i2}\phi_2$ for some $a_{i1}$, $a_{i2}$, the matrix $\Sigma$ can be obtained from the matrix K with ijth entry $<\xi_i,\xi_j>$ by the addition of some rank 2 matrix which is not important to our problem. (The $a_{ij}$ depend on the definition of $<\phi_\mu,\phi_\nu>$, $\mu,\nu = 1,2$, which is irrelevant to the estimate being studied.) Furthermore, if $\gamma_1 \geq \gamma_2 \geq \ldots \geq \gamma_n$ are the eigenvalues of K, and $\delta_1 \geq \delta_2 \geq \ldots \geq \delta_{n-2}$ are the eigenvalues of QKQ' then $\gamma_{\nu-2} \leq \delta_\nu \leq \gamma_\nu$. Now as part of the calculations for Sections 3 and 4 we have computed $b_1^2 \geq \ldots b_{n-2}^2$, which are the eigenvalues of $Q\tilde{K}Q'$ $\tilde{K}$ being the n×n matrix with ijth entry $<P_N\xi_i,P_N\xi_j>$. The $b_\nu^2$'s satisfy $b_\nu^2 \leq \gamma_\nu$, i=1,2,...,n-2 and the number of non zero $b_\nu^2$'s cannot be bigger than the dimension of the range of $P_N$. In the limit as N→∞, $b_\nu^2 \to \gamma_\nu$. If N is too small, it, too, can act as a "smoothing parameter".

Figure 5.1 gives a plot of the first 68 $b_\nu^2$'s on a log-log plot, with n = 80, N = 80. (The vertical unit is arbitrary and depends on the units in which r is carried in the computer. It is reasonable to choose these units so that $b_1^2 \approx 1$). For comparison an arrow marks $n\hat{\lambda} = 80 \times 10^{-5}$. Recall that the eigenvalues of $\tilde{A}(\lambda)$ are $(1,1,b_1^2/(b_1^2+n\lambda),\ldots, b_{n-2}^2/(b_{n-2}+n\lambda))$. $\tilde{A}(\lambda)$ plays the role of the influence matrix $X(X'X+n\lambda I)X'$ in the regression problem when a ridge estimate is used for $\beta$.
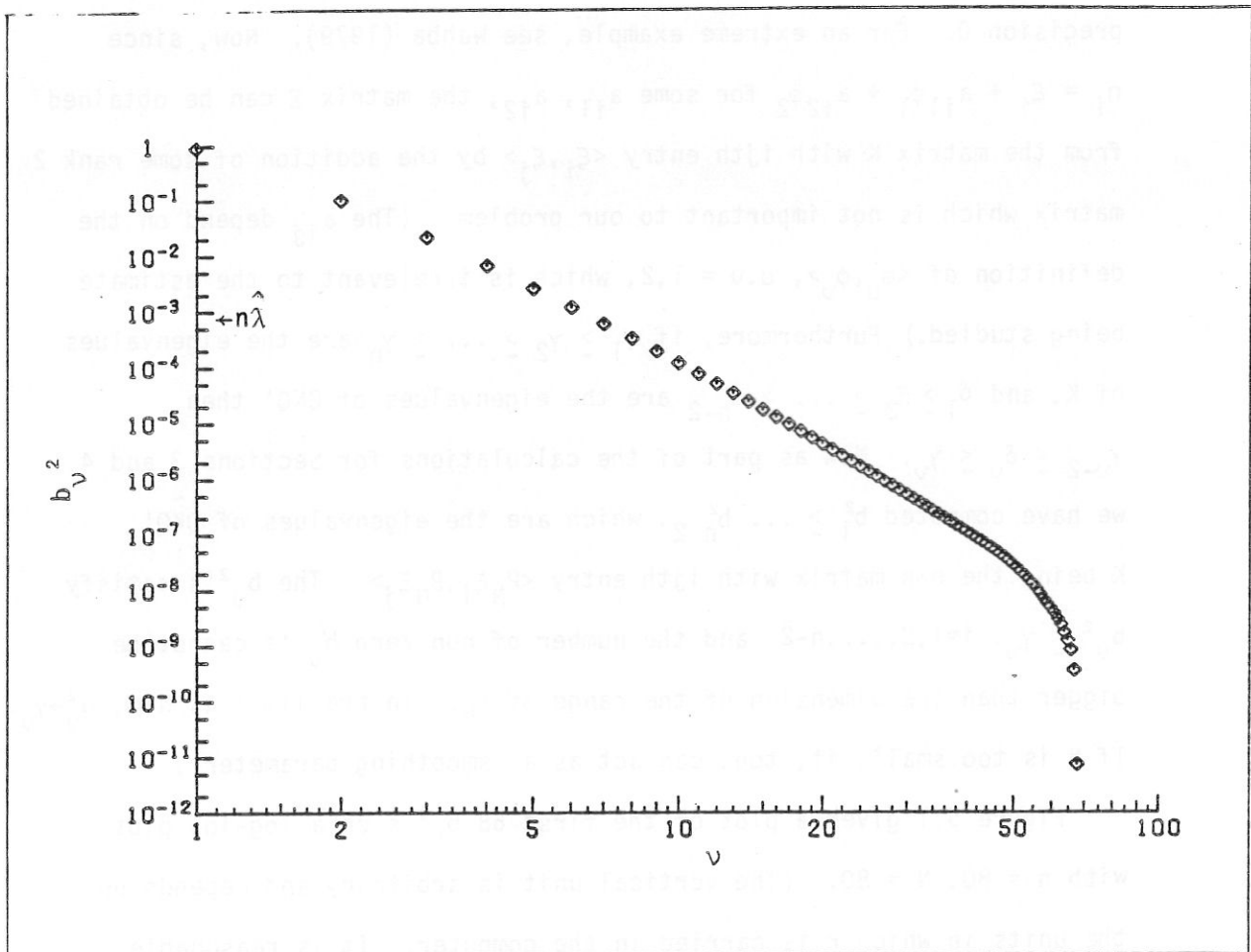
Figure 5.1  The eigenvalues of QKQ'; n = 80, N = 80

Based on trying several values of n and N, it is our belief that at least the first 30 or 40 $b_\nu^2$'s approximate the $\delta_\nu$'s very well, and that increasing N would have no appreciable effect on the resulting estimate $f_{\hat\lambda}$. If n is increased, our unpublished plots as well as recent analytical work suggests that the slope of the (major part of) the eigensequence log-log plot will tend to a limit. (See, e.g. Utreras (1981), Wahba (1977).) Note that $b_{40}^2/b_1^2$ is already down to $10^{-7}$. We conclude that increasing n (with N>n) much past 80 would not change $f_{\hat\lambda}$, certainly not to plot accuracy, and that we have thus succeeded in choosing n and N so they are not acting as smoothing parameters.

Eigensequence plots can provide insight about practical limits on the amount of information concerning $f_3^\epsilon$, in the data, and we suggest that these plots be routinely examined in problems of this sort. It is seen that with $\hat\lambda = 10^{-5}$, the eigenvalues $b_\nu^2/(b_\nu^2+n\lambda)$ of the influence matrix $A(\lambda)$ have decreased to .5 by about the 8th eigenvalue ($\nu=6$),

## 7. RELATED ESTIMATES AND THEIR SMOOTHING PARAMETERS

Another approach to the approximate numerical calculation of the minimizer of (2.2) in $H$ is to minimize (2.2), in a convenient approximating subspace $S_N = \text{span}\{B_\ell\}$, say.

Then one finds $h_\lambda$ of the form

$$h_\lambda = \sum_\ell \theta_\ell B_\ell$$

to minimize (2.2). In the problems studied here a space of cubic B-splines (see deBoor (1978)) would be appropriate. If the basis functions have compact support, this would be considered to be a "finite element" method.

Given $\varepsilon = s_0 < s_1 < \ldots s_N = R$, let $S_N h$ be that function in $H$ which minimizes $J(h)$ subject to $(S_N h)(s_\ell) = h(s_\ell)$, $\ell = 0,1,\ldots,N$, and $(S_N h)'(s_0) = h'(s_0), (S_N g)'(s_N) = h'(s_N)$. $S_N h$ is the cubic spline interpolating to $h$ at $s_0, s_1, \ldots, s_N$, and to $h'$ at $s_0$ and $s_N$. Let $S_N$ be a set of $N+3$ cubic B-splines whose span is the range of $S_N$. (See, e.g. deBoor (1978) Chapter IX.) Then the minimizer of the exact expression

$$\frac{1}{n} \sum_{i=1}^{n} (<\eta_i, h> - z_i)^2 + \lambda J(h) \tag{7.1}$$

in the approximating subspace $S_N$, is the same as the minimizer of the approximate expression

$$\frac{1}{n} \sum_{i=1}^{n} (<\eta_i, S_N h> - z_i)^2 + \lambda J(h) \tag{7.2}$$

in $H$. This can be shown without difficulty by writing $h = S_N h + (I - S_N)h = \sum_\ell \theta_\ell B_\ell + (I - S_N)h$ for some $\{\theta_\ell\}$ and using the property of the cubic spline interpolant $\int_\varepsilon^R (S_N h)''((I - S_N)h)'' = 0$, to obtain

$$J(h) = J(S_N h) + J((I - S_N)h).$$

It can then be shown that the minimizer of (7.2) must be in $S_N$. Upon observing that $S_N h = h$ for any $h$ in $S_N$ (spline interpolation is idempotent) it follows that problems (7.1) and (7.2) are the same. For comparison with (2.11) we can write (7.2) as

$$\frac{1}{n} \sum_{i=1}^{n} (\tilde{\tilde{L}}_i h - z_i)^2 + \lambda J(h) \tag{7.3}$$

where $\tilde{\tilde{L}}_i h = \langle S_N^* \eta_i, h \rangle$, $S_N^*$ being the adjoint operator to $S_N$. The cross validation function $\tilde{V}(\lambda)$ for the problem (7.3) can also be readily obtained.

The minimizer of the exact problem (7.1) in some space $S_N$ of B-splines is the "hybrid estimate" proposed by Wahba (1980) and mentioned by Mendelsohn and Rice (1983). If the dimension of $S_N$ is chosen large then this "hybrid estimate" will numerically be a good approximation to the original cross validated spline estimate (minimizer of (2.2) in $H$). On the other hand, if $S_N$ is relatively small, then $N$ will act as a smoothing parameter. Thus there will be a pair of smoothing parameters $(\lambda, N)$, which, in principle, could be chosen objectively by GCV. Mendelsohn and Rice (1983) solved the problem mentioned in (1.7) by using the minimizer of

$$\frac{1}{n} \sum_{i=1}^{n} (L_i h - z_i)^2$$

in $S_N$. In their work $n$ was very large, and $N$ was the only smoothing parameter.

When $N$ is the only smoothing parameter, the optimal (IMSE) value of $N$ grows very slowly with $n$. (For certain regression problems $N = 0(n^{1/5})$, see Agarwal and Studden (1980)). When $n$ is very large and the data is nearly exact, then one can sometimes profitably use $N$ as the sole smoothing parameter,

since the optimal N will be large enough so that recoverable structure in the solution will not be lost.  (An "N only" estimate is the easiest to compute.)  In Mendelsohn and Rice's problem, n was several hundred and the data could be considered extremely "exact" since $10^5$ observations were in the n bins.  They found an N of 12 subjectively.  In our problem with much "noisier" data we conjecture that the optimal N in an N only estimate would result in N of more like 3-6, and in general the estimate would not show the peak resolution that is evident in Figures 3 and 4 unless the true solution was actually in $S_N$.  Efficient numerical methods for the hybrid estiamte for problems with very large N (as might occur in image processing, for example) can be found in Bates and Wahba (1982).

We see now that there are actually three possible smoothing parameters, $\lambda$, n and N.  In the "matched quadrature" method, it is natural to have N > n and the computing load is sensitive to n and insensitive to N.  In the "hybrid method" it is natural to take N < n and the computing load will be sensitive to N and insensitive to n.  In the "matched quadrature" method, one could easily use n and $\lambda$ as joint smoothing parameters and in the "hybrid method" one could easily use N and $\lambda$ as joint smoothing parameters.  (There may, however, be a region in $(\lambda,n)$ or $(\lambda,N)$ space where decreasing both parameters simultaneously will have little effect on the IMSE.)  In the problem at hand, where the data is very noisy (because the sample size is small) and the problem is somewhat ill posed, we believe that one can do a better job of recovering structure in the solution if one lets $\lambda$ do all, or most of the smoothing, and one chooses n and N just large enough so that they are not doing appreciable

smoothing.  When there is a very large amount of information in the data,
using n and/or N to do (some of) the smoothing, can be very cost effective.

## Appendix A

Formula for $\xi_i(r)$ and $\tau_{i\nu}$

$\xi_i(r) = L_i(Q_1(\cdot,r)) = \psi(P_{i-1},r) - \psi(P_i,r)$, where

$$\psi(P,t) = \frac{(t-\epsilon)^2}{2} J_1(\epsilon,P,P,R) - \frac{(t-\epsilon)^3}{2} I_0(P,P,R), \quad \epsilon \le t \le P$$

$$= \frac{(t-\epsilon)}{2} J_2(\epsilon,P,P,t) - I_3(\epsilon,P,P,t)$$

$$+ \frac{(t-\epsilon)}{2} J_1(\epsilon,P,t,R) - \frac{(t-\epsilon)^3}{6} I_0(P,P,R) \quad P \le t \le R$$

where

$$I_k(x,a,b) = \int_a^b u^k \sqrt{u^2-x^2}\,du \qquad x \le a, \; k = 0,1,2,3$$

$$J_k(\epsilon,x,a,b) = \sum_{i=0}^k \binom{k}{i} \epsilon^{k-i} I_i(x,a,b) \qquad k = 0,1,2,3.$$

The definite integrals $I_k$ have closed form analytic representation, see Selby (1979) Formulae No.'s 156, 167, 168 and 170, p. 425.

$$\tau_{i1} = I_0(P_{i-1},P_{i-1},R) - I_0(P_i,P_i,R)$$

$$\tau_{i2} = J_1(P_{i-1},P_{i-1},R) - J_1(P_i,P_i,R)$$

## Appendix B

## Computation of $P_N \xi_i$ and $\tilde{k}_{ij}$

Since $\xi_i(0) = \xi_i'(0) = 0$, and $Q_1$ is the reproducing kernel for the subspace of $H$ satisfying these boundary conditions, we must have

$$P_N \xi_i = \sum_{k=1}^{N} \alpha_{ik} Q(\cdot, s_k),$$

for some $\alpha_i = (\alpha_{i1}, \alpha_{i2}, \ldots, \alpha_{iN})$. The $\alpha_{ij}$ are chosen so that the interpolation conditions are satisfied, that is

$$\tilde{Q}\alpha_i = \begin{pmatrix} \xi_i(s_1) \\ \vdots \\ \xi_i(s_N) \end{pmatrix}$$

where $\tilde{Q}$ is the N×N matrix with ijth entry $Q(s_i, s_j)$. Since $\tilde{Q}$ is positive definite, the $\alpha_i$ can be efficiently computed via a Cholesky factorization of $\tilde{Q}$ (see Dongarra et al. (1978), Chapter 3).

Now $\tilde{k}_{ij} = \langle P_N \xi_i, P_N \xi_j \rangle = \sum_k \alpha_{ik} \sum_\ell \alpha_{j\ell} \langle Q(\cdot, s_k), Q(\cdot, s_\ell) \rangle$

$$= \alpha_i' \tilde{Q} \alpha_j.$$

# REFERENCES

Agarwal, G. and Studden, W. (1980). Asymptotic integrated mean square error using least squares and bias minimizing splines. Ann. Statist. 8, 6, 1307-1325.

Anderssen, R.S. and Jakeman, A.S. (1975). Abel type integral equations in stereology, II. Computational methods of solution and the random spheres approximation. J. Microscopy 105, 135-153.

Bates, D. and Wahba, G. (1982). Computational methods for generalized cross-validation with large data sets. In "Treatment of Integral Equations by Numerical Methods". C.T.H. Baker and G.F. Miller, eds. Academic Press, London, 283-296.

Cox, D.R. (1970). Analysis of Binary Data, Chapman and Hall, London.

Cox, Dennis (1983). Draft manuscript.

Chover, J., King, J. (1981). Personal communication.

Craven, P. and Wahba, G. (1979). Smoothing noisy data with spline functions: estimating the correct degree of smoothing by the methods of generalized cross validation. Numer. Math. 31, 377-403.

Crump, J.G., and Seinfeld, J.H. (1982). A new algorithm for inversion of aerosol size distribution data. Aerosol Science and Technology 1: 15-34.

deBoor, C. (1978). A Practical Guide to Splines. Springer-Verlag, New York.

Diaconis, P., and Efron, B. (1983). Computer-intensive methods in statistics, Scientific American 248, 5, (May) 116-130.

Dongarra, J.J., Bunch, J.R., Moler, C.B., and Stewart, G.W., (1979). LINPACK users guide. Society for Industrial and Applied Mathematics, Philadelphia.

Gasser, T. and Muller, M. (1979). Kernel estimation of regression functions. In "Smoothing Techniques for Curve Estimation", T. Gasser and M. Rosenblatt, eds., Lecture Notes in Mathematics No. 757. Springer-Verlag.

Keiding, N., Jensen, S.T. and Ranek, L. (1972). Maximum likelihood estimation of the size distribution of liver cell nuclei from the observed distribution in a plane section. Biometrics 28, 813-829.

Kimeldorf, G., and Wahba, G. (1971). Some results on Tchebycheffian spline functions. J. Math. Anal. Applic. 33, 1, 82-95.

Koen, H., Pugh, T, and Goldfarb, S, (1983), Hepatocarcinogenesis in the Mouse Combined Morphologic-Stereologic Studies. To appear Am. J. Path, July, 1983, p. 89-100.

Kuk, A.Y.C. (1982). A mixing distribution approach to estimating particle size distributions. Stanford University Statistics Department Technical Report No. 328.

Lukas, M. (1981). Regularization of linear operator equations, Thesis, Australian National University.

Mendelsohn, J. and Rice, J. (1982). Deconvolution of microfluorometric histograms with B-splines. J. Am. Stat. Assoc. 77, 380, 748-753.

Merz, P. (1980). Determination of adsorption energy distribution by regularization and a characterization of certain adsorption isotherms. J. Comput. Physics, 38, 64-85.

Minerbo, G.N., Levy, M.E. (1969). Inversion of Abel's integral equation by means of orthogonal polynomials. SIAM J. Numer. Anal. 6; 598-616.

Nicholson, W.L. (1976). Estimation of linear functions by maximum likelihood. J. Microscopy 113, 113-239.

Nicholson, W.L. (1970). Estimation of linear properties of particle size distributions. Biometrika, 57, 273-297.

Nicholson, W.L., and Merck, . (1969). Unfolding particle size distributions. Technometrics II, 707-720.

Nychka, D. (1983). The analysis of some penalized likelihood estimation schemes, Thesis, Statistics Department, University of Wisconsin-Madison, to appear.

Selby, S., ed. (1979). CRC Standard Mathematical Tables, 21st Edition. The Central Rubber Co., Cleveland.

Silverman, B. (1983). On the estimation of a probability density function by the maximum penalized likelihood method. Ann. Statist. 10, 3, 795-810.

Smith, B.T., Boyle, J.M., Dongarra, J.J., Garbow, B.S., Ikebe, Y., Klema, V.C. and Moler, C.B. (1976). Matrix Eigensystem Routines-- EISPACK Guide. Lecture notes in Computer Science, Springer-Verlag.

Tallis, G.M. (1970). Estimating the distribution of spherical and elliptical bodies in conglomerates from plate sections. Biometrics 26, 87-103.

Utreras, F. (1981). Optimal smoothing of noisy data using spline functions, SIAM J. Sci. Stat. Comput., 2, 3, 349-362.

Villalobos, M. (1983). Thesis, to appear.

Villalobos, M., and Wahba, G. (1983). Multivariate thin plate spline estimates for the posterior probabilities in the classification problem, to appear, Commun. Stat. B.

Wahba, G. (1975). Interpolating spline methods for density estimation. I. Equispaced knots. Ann. Statist. 3, 30-48.

Wahba, G. (1977).  Practical approximate solutions to linear operator equations when the data are noisy.  SIAM J. Numer. Anal., 14, 4.

Wahba, G. (1978).  Improper priors, spline smoothing and the problem of guarding against model errors in regression.  J. Roy. Stat. Soc. Ser. B., 40, 3.

Wahba, G. (1979).  Smoothing and ill posed problems, in "Solution Methods for Integral Equations with Applications", Michael Golberg, ed., 183-194, Plenum Press, (1979).

Wahba, G. (1980).  "Ill posed problems:  Numerical and statistical methods for mildly, moderately and severely ill posed problems with noisy data.  TR No. 595, University of Wisconsin-Madison, Statistics Department.

Wahba, G. (1982a).  Constrained regularization for ill posed linear operator equations with applications in meteorology and medine in "Statistical Design Theory and Related Topics: III, Vol. 2, S.S. Gupta and J.O. Berger, eds., Academic Press.

Wahba, G. (1982b).  Cross validated spline methods for direct and indirect sensing experiments.  University of Wisconsin Statistics Department TR #694, to appear in "Proceedings of the Signal Processing in the Ocean Environment Workshop", E.J. Wegman, ed., Marcel-Dekker, Inc.

Wahba, G. and Wendelberger, J. (1980).  "Some new mathematical methods for variational objective analysis using splines and cross-validation.  Monthly Weather Review 108, 1122-1143.

Watson, G.S. (1971).  Estimating functionals of particle size distributions.  Biometrika 58, 483.

Wendelberger, J. (1981).  The computation of Laplacian smoothing splines with examples.  University of Wisconsin Statistics Department Technical Report No. 648.

Wicksell, D.S. (1925).  The corpuscle problem, Part I.  Biometrika 17, 87-97.