1.	AD-	A056 70	A C MAR	ATE UNI DENSITY 78 E 60	V OF NE	W YORK	AT BUF	FALO AN	HERST	STATIS ROBUST D	TICAE ESTIMA AAG29-7	TC F/0 TION.(0 76-6-02	G 12/1 U) 39 L	
		0F   A056707					A second						1000 100 1000 1	
				2 5 - 2 - 2 - 2 - 2 - 2 - 2 - 2 - 2 - 2 -						$\begin{array}{c} \begin{array}{c} & & & & & \\ & & & & \\ & & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & $	<ul> <li>Bernsteinen Sternen Ster</li></ul>	2 - (1) - (2) - (2)		
	1 1 1 1 1	ан 				$\label{eq:second} \begin{array}{c} & & & \\ & & & & \\ & & & \\ & & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ &$		<ul> <li>A manufacture of the second se</li></ul>		70 - 1 - 2 - 2 - 2 - 2 - 2 - 2 - 2 - 2 - 2			Önge Føge	Gede Gede
	lelele Idda	adada Adada		Freedom?	END DATE FILMED 9 = 78						¢			
		2	Harry America						2					
	1													. /



State University of New York at Buffalo **Department of Computer Science** AD A 0 5670 ARO 13845.5-N A DENSITY-QUANTILE FUNCTION PERSPECTIVE ON ROBUST ESTIMATION by 2000000 10 Emanuel/Parzen JUL 26 1978 Statistical Science Division troi T 1511 State University of New York at Buffalo An Invited Paper to be Published in the Proceedings of the ARO Workshop in Robustness in Statistics AU NO. GRANT TECHNICAL REPORT NO. ARO-4 STATISTICAL SCIENCE DIVISION REPORT NO. 60 March 1978 \*Research supported by the Army Research Office (Grant/DA AG29-76-0239). Approved for public release; distribution unlimited. The findings in this report are not to be construed as an official Department of the Army position. unless so designated by other authorized documents. 78 07 21 030 409 511



# A DENSITY-QUANTILE FUNCTION PERSPECTIVE ON ROBUST ESTIMATION\*

by

Emanuel Parzen Institute of Statistics Texas A & M University

BY DISTRIBUTION/AVAILABILITY CODES Dist. AVAIL and/or SPECIAL	MTIS DDC UNANNOUNI JUSTIFICAT	White Section Buff Section ED
Dist. AVAIL. and/or SPECIAL	BY. DISTRIBUTI	ON/AVAILABILITY CODES
	the state of the s	

## 1. Introduction

Parametric statistical inference may be said to be concerned with statistical inference of idealized parameters from ideal data. Huber (1977), p. 1, writes: "The traditional approach to theoretical statistics was and is to optimize at an idealized parametric model."

Robust statistical inference may be said to be concerned with statistical inference of idealized parameters from semi-ideal data (by the use of methods which are insensitive against small deviations from the ideal assumptions). Huber (1977), p. 3, writes: the robust approach to theoretical statistics assumes "an idealized parametric model, but in addition one would like to make sure that methods work well not only at the model itself, but also in a neighborhood of it."

Exploratory data analysis may be said to be concerned with statistical inference from non-ideal data(often by seeking re-expressions (transformations) of the data that will make it more ideal). Exploratory data analysis helps pose the well-posed statistical questions to which classical parametric statistics provides answers.

\*Research supported by the Army Research Office (Grant DA AG29-76-0239).

This paper provides an overview to a new general approach to statistical data analysis and parameter estimation which could be called the <u>quantile</u> function approach. The aims of <u>descriptive statistics</u> (to graphically summarize and display the data) are obtained by Quantile-Box plots of the sample quantile function. The aims of <u>"goodness of fit"</u> are obtained by fitting smooth quantile functions to the sample quantile function. The aims of <u>parameter estimation</u>, especially robust estimation of location and scale parameters, are attained by regression analysis of the sample quantile function. (The goal of a statistician in analyzing a batch of data  $X_1, \ldots, X_n$  should be both "estimation of parameters" and "goodness of fit". By "goodness of fit" is meant fitting of the observed sample probabilities by a smooth probability law.)

Quantile functions are defined in section 2. Window estimators of location and scale parameters are defined in section 3; their equivalence to L-estimators is discussed in section 4. A conjectured expression is given in section 5 for the asymptotic variance of window estimators. New approaches being developed for non-parametric probability law modeling are mentioned in section 6; quantile box-plots are introduced in section 7. Section 8 discusses location and scale parameter estimation using trimmed samples. Robust regression is the subject of section 9. A new definition of statistics is proposed in section 10.

To carry out in practice robust estimation of location parameters this paper proposes computing <u>means which adapt to the ends</u> (by "ends"

78 07 21 030

one means the tail character of the distribution of the data). Three such methods are given in the paper:

(1) Iteratively reweighted estimators with weight function  $w(x) = (1 + \frac{1}{m}x^2)^{-1}$  for suitable choices of m (section 3);

(2) Maximum likelihood estimation omitting extreme order statistics where the percentage of values omitted is determined from the goodness of fit of the corresponding smooth quantile functions (section 8);

(3) Adaptive L-estimation of location and scale parameters using autoregressive estimators of density-quantile functions (section 8).

A fourth method of robust location and scale parameter estimation is:

(4) Quantile box-plot diagnostics which indicate that mid-summaries and mid-scales are equal enough to provide naive estimators of location and scale (section 7).

# 2. Quantile Function

The <u>quantile</u> function Q(u),  $0 \le u \le 1$  of a random variable X is the <u>inverse</u> of its distribution function  $F(x) = P(X \le x)$ . The precise definition of Q is:

$$Q(u) = F^{-1}(u) = \inf \{x: F(x) \ge u\}$$
.

Given a sample  $X_1, \ldots, X_n$ , we denote the sample distribution function by  $\tilde{F}(x)$ ,  $-\infty < x < \infty$ ; it is defined by

$$F(x) = fraction of X_1, \dots, X_n \le x$$
.

The Sample Quantile Function

$$\widetilde{Q}(u) = \widetilde{F}^{-1}(u) = \inf \{x: \widetilde{F}(x) \ge u\}$$

can be computed explicitly in terms of the order statistics  $X_{(1)} < X_{(2)} < \dots$ <  $X_{(n)}$  (which are the values in the sample arranged in increasing order):

$$\widetilde{Q}(u) = X_{(j)}, \quad \frac{j-1}{n} < u \leq \frac{j}{n}$$

The foregoing definition of  $\tilde{Q}(u)$  is a piecewise constant function. It is more convenient to define  $\tilde{Q}(u)$  as a <u>piecewise linear</u> function. Divide the unit interval into 2n subintervals. For u = (2j - 1)/2n define

$$\tilde{Q}(\frac{2j-1}{2n}) = X_{(j)}, \quad j = 1, 2, \dots n$$

For u in  $\frac{2j-1}{2n} \le u \le \frac{2j+1}{n}$ , j = 1, 2, ..., n-1,

define  $\widetilde{Q}(u)$  by linear interpolation; thus for u in this interval

$$\widetilde{Q}(u) = n(u - \frac{2j-1}{2n}) X_{(j+1)} + n(\frac{2j+1}{2n} - u) X_{(j)}$$

In particular

$$\tilde{Q} \left(\frac{j}{2n}\right) = \frac{1}{2} X_{(j+1)} + \frac{1}{2} X_{(j)}$$

The population median is Q(0.5). The sample median is  $\tilde{Q}(0.5)$ . Our definition of  $\tilde{Q}(u)$  has the merit that  $\tilde{Q}(0.5)$  is the usual definition of the sample median:

$$\widetilde{Q}(0.5) = X_{(m+1)}$$
 if  $n = 2m + 1$  is odd,  
=  $\frac{1}{2} (X_{(m)} + X_{(m+1)})$  if  $n = 2m$  is even.

The asymtotic distribution of Q(u) satisfies:  $\sqrt{n} fQ(u) \{Q(u) - Q(u)\}$ is asymptotically normal, with mean 0 and variance u(1 - u), where fQ(u) denotes the probability density function f(x) = F'(x)evaluated at x = Q(u); in symbols,

$$fQ(u) = f(Q(u))$$

We call fQ(u) the density-quantile function.

Estimating the fQ-function is of interest for two reasons: as a way of estimating (1) the true probability density function f(x), and (2) approximate confidence intervals for Q(u) and especially for the true median Q(0.5), since

$$\tilde{Q}(0.5) \pm \{\sqrt{n} fQ(0.5)\}^{-1}$$

is an approximate 95% confidence interval for the median Q(0.5).

We call q(u) = Q'(u) the quantile-density function. The identity

FQ(u) = u

implies the reciprocal relationship

fQ(u) = 1.

Thus we may write (using  $\doteq$  to denote approximate equality)

$$\frac{1}{fQ(0.5)} = q(0.5) \doteq \frac{Q(0.75) - Q(0.25)}{0.75 - 0.25} = 2\{Q(0.75) - Q(0.25)\}$$

7

We define, for  $0 \le p \le 0.5$ ,

$$R(p) = Q(1-p) - (p)$$

to be the p-range, and

$$\widetilde{R}(p) = \widetilde{Q}(1-p) - Q(p)$$

to be the sample p-range. When p = 0.25, we call Q(0.75) and Q(0.25) the quartiles,

$$R(0.25) = Q(0.75) - Q(0.25)$$

the quartile-range, and

$$\tilde{R}(0.25) = \tilde{Q}(0.75) - \tilde{Q}(0.25)$$

the sample quartile-range.

One can conclude that the median Q(0.5) has a non-parametric estimator given by  $\widetilde{Q}(0.5)$ , and an approximate 95% confidence interval given by

$$\tilde{Q}(0.5) \pm 2\tilde{R}(0.25)/\sqrt{n}$$

A use of a confidence interval of this kind for the median is discussed by McGill, Tukey, and Larsen (1978).

The aim of the foregoing discussion is to introduce the quantile function and illustrate how it is traditionally used to provide non-parametric measures of location (such as the median) and scale (such as the quartile range). Our aim is to use quantile functions to detect and describe ideal and non-ideal statistical models for data.

# 3. Location and Scale Estimation by Window Estimators

One of the points which this paper would like to make is that measures of location and scale of a data sample are interpretable only if they are <u>probability based</u>, in the sense that they are estimators of characteristics of the true quantile function of the random variable X.

We use  $\mu$  and  $\sigma$  to denote measures of location and scale respectively. When  $\mu$  and  $\sigma$  represent median and inter-quartile range,  $\mu = Q(0.5)$ and  $\sigma = Q(0.75) - Q(0.25)$ . When  $\mu$  and  $\sigma^2$  represent mean and variance, they can be expressed in terms of Q by

$$\mu = \int_0^1 Q(u) du$$
,  $\sigma^2 = \int_0^1 \{Q(u) - \mu\}^2 du$ 

These formulas follow immediately from the basic fact that X is identically distributed as Q(U) where U is uniformly distributed on the interval [0,1].

When  $\mu$  and  $\sigma^2$  represent mean and variance, fully non-parametric estimators of  $\mu$  and  $\sigma^2$  are

$$\tilde{\mu} = \int_0^1 \tilde{Q}(u) \, du$$
,  $\tilde{\sigma}^2 = \int_0^1 {\{\tilde{Q}(u) - \tilde{\mu}\}}^2 \, du$ 

which are essentially the sample mean and the sample variance.

To efficiently estimate location and scale parameters  $\mu$  and  $\sigma$ , it is customary to start with a model for the probability density function f(x) of the form

$$f(\mathbf{x}) = \frac{1}{\sigma} f_0 \left(\frac{\mathbf{x} - \mu}{\sigma}\right) \tag{*}$$

where  $f_0(x)$  is a known probability density function. Define  $L(\mu, \sigma)$  to be (1/n) times the log - likelihood of the sample  $X_1, \ldots, X_n$ ; it is given by

$$L(\mu,\sigma) = -\log \sigma + \frac{1}{n} \sum_{i=1}^{n} \log f_0(\frac{X_i - \mu}{\sigma})$$

One can express likelihood in terms of quantile functions:

$$L(\mu,\sigma) = -\log \sigma + \int_0^1 \log f_0(\frac{\tilde{\Omega}(u) - \mu}{\sigma}) du$$
.

The model (\*) leads to a very simple formula for the true quantile function Q(u) of the data:

$$Q(u) = \mu + \sigma Q_0(u)$$

where  $Q_0(u)$  is a known quantile function corresponding to  $f_0(x)$ . For ease of writing we introduce the notation

$$\tilde{Q}_0(u) = \frac{Q(u) - \mu}{\sigma}$$

The maximum likelihood estimators  $\mu$  and  $\sigma$  satisfy the log likelihood-derivative equations:

$$\frac{\partial}{\partial \mu} L(\mu, \sigma) = 0$$
,  $\frac{\partial}{\partial \sigma} L(\mu, \sigma) = 0$ 

To compactly write formulas for these derivatives, define

$$\psi(\mathbf{x}) = \frac{-f_0'(\mathbf{x})}{f_0(\mathbf{x})} = -\frac{\partial}{\partial(\mathbf{x})} \log f_0(\mathbf{x}) .$$

 $\mathbf{w}(\mathbf{x}) = \frac{1}{\mathbf{x}} \psi(\mathbf{x})$ 

$$\frac{\partial}{\partial \mu} L(\mu, \sigma) = \frac{1}{\sigma} \int_0^1 \psi(\tilde{\Omega}_0(u)) du$$
$$= \frac{1}{\sigma^2} \int_0^1 \{\tilde{\Omega}(u) - \mu\} w(\tilde{\Omega}_0(u)) du$$
$$\frac{\partial}{\partial \sigma} L(\mu, \sigma) = -\frac{1}{\sigma} + \frac{1}{\sigma^2} \int_0^1 \psi(\tilde{\Omega}_0(u)) \{\tilde{\Omega}(u) - \mu\} du$$
$$= -\frac{1}{\sigma} + \frac{1}{\sigma^3} \int_0^1 w(\tilde{\Omega}(u)) \{\tilde{\Omega}(u) - \mu\}^2 du$$

In the normal case,  $\psi(x) = x$ , w(x) = 1 and  $\mu$  and  $\sigma^2$  are equal to the sample mean and variance respectively.

To obtain estimators  $\mu$  and  $\sigma$  without specifying  $f_0(x)$ , one introduces the concept of <u>iteratively reweighted estimators</u> of  $\mu$  and  $\sigma^2$ . Given estimators  $\mu^*$  and  $\sigma^*$  define

$$\tilde{Q}_0^*(u) = \frac{\tilde{Q}(u) - \mu^*}{\sigma^*}$$

Then as "approximate" solutions of the log-likelihood derivative equations, one studies the estimators defined by

$$\hat{\mu} = \frac{\int_{0}^{1} \tilde{Q}(u) w(\tilde{Q}_{0}^{*}(u)) du}{\int_{0}^{1} w(\tilde{Q}_{0}^{*}(u)) du}$$
$$\hat{\sigma}^{2} = \int_{0}^{1} \{\tilde{Q}(u) - \hat{\mu}\}^{2} w(\tilde{Q}_{0}^{*}(u)) du .$$

These formulas for  $\mu$  and  $\sigma$  reduce to the sample mean and variance when one chooses  $w(x) \equiv 1$ .

Since we are concerned with forming estimators of location and scale which are satisfactory for long-tailed distributions it is natural to choose weight functions w(x) corresponding to Students' t-distribution with m degrees of freedom,

$$f_0(x) = \frac{1}{\sqrt{m\pi}} \frac{\Gamma(\frac{m+1}{2})}{\Gamma(\frac{m}{2})} \left(1 + \frac{x^2}{m}\right)^{-(m+1)/2}$$

for which

$$w(x) = -\frac{1}{x} (\log f_0(x))' = \frac{m+1}{m} \frac{1}{1+\frac{1}{m}x^2}$$

We call this weight function a window, and we call  $\mu$  and  $\sigma$  window estimators.

To completely specify the window, one must specify a value for m (which we could call the "trimming width" of the window). The more normal the data is believed to be, the larger should m be chosen (say, m = 25). The more Cauchy-distributed the data is believed to be, the closer to 1 should m be chosen (say, m = 4). In practice, one might try both values of m, and compare the results. The constant m could also be estimated adaptively to yield "self-tuning" robust estimators of location and scale. A window recommended by Tukey (see Mosteller and Tukey (1977), p. 205) is the bisquare window:

$$w_{Bisquare}(x) = (1 - (\frac{x}{c})^2)_+^2$$

where c is a suitably chosen constant. Tukey recommends that c be taken to be 6 or 4 when x is measured in units of  $\sigma$ . It seems likely that the choice of c should reflect one's beliefs about the long-tailed character of the data.

# 4. Weight Functions of L-Estimators

An L-estimator  $\mu$  of a location parameter is a linear combination of order statistics  $X_{(1)} < \ldots < X_{(n)}$ , which we write in the form

$$\hat{\mu} = \int_0^1 \tilde{Q}(u) W(u) du$$

for suitable weight function W(u). Asymptotically efficient L-estimators of  $\mu$  and  $\sigma$  in the model  $Q(u) = \mu + \sigma Q_0(u)$ , when  $f_0$  is a symmetric density, are given by [see Parzen (1978), and summary in section 8]

$$\hat{\mu} = \int_0^1 \widetilde{Q}(u) \ W_{\mu}(u) \ du + \int_0^1 \ W_{\mu}(u) \ du$$
$$\hat{\sigma} = \int_0^1 \widetilde{Q}(u) \ W_{\sigma}(u) \ du + \int_0^1 \ W_{\sigma}(u) \ du$$
where
$$W_{\mu}(u) = f_0 Q_0(u) \ J_0'(u) = \frac{J_0'(u)}{Q_0'(u)} \ and$$

$$W_{\sigma}(u) = J_{0}(u) + Q_{0}(u) W_{\mu}(u)$$
.

 $f_0 Q_0(u)$  is the density-quantile function corresponding to  $Q_0$ , and  $J_0(u)$ is its score function defined by

$$J_{0}(u) = -(f_{0}Q_{0})'(u) = \frac{-f_{0}'Q_{0}(u)}{f_{0}Q_{0}(u)} = \psi(Q_{0}(u))$$

An L-estimator forms a weighted average of order statistics in which the weights depend on the <u>ranks</u> u. It is of interest to express the weights as a function of  $Q_0(u)$ , which is the <u>size</u> of the order statistics. One can derive such formulas starting from the general representation, given by Parzen (1978),

$$f_0 Q_0(u) \sim (1-u)^{\alpha}$$
,  $Q_0(u) \sim (1-u)^{-(\alpha-1)}$ 

where  $\alpha$ , called the tail exponent, is assumed to satisfy  $\alpha > 1$  (indicative of long tailed distributions). We write

$$\begin{split} \mathbf{W}_{\mu} &= \mathbf{J}_{0}^{'} \mathbf{f}_{0} \mathbf{Q}_{0} &= -(\mathbf{f}_{0} \mathbf{Q}_{0})^{2} [(\log \mathbf{f}_{0} \mathbf{Q}_{0})^{"} + \{(\log \mathbf{f}_{0} \mathbf{Q}_{0})^{'}\}^{2}] \\ \mathbf{W}_{\sigma} &= \mathbf{J}_{0} + \mathbf{Q}_{0} \mathbf{W} = -(\mathbf{f}_{0} \mathbf{Q}_{0})(\log \mathbf{f}_{0} \mathbf{Q}_{0})^{'} + \mathbf{Q} \mathbf{W}_{\mu} \end{split}$$

Therefore

$$W_{\mu}(u) \sim (1-u)^{2(\alpha-1)} \alpha(1-\alpha) \sim \frac{1}{Q_{0}^{2}(u)} \alpha(1-\alpha)$$
$$W_{\sigma}(u) \sim (1-u)^{\alpha-1} \alpha(2-\alpha) \sim \frac{1}{Q_{0}(u)} \alpha(2-\alpha) .$$

The main conclusion we desire to point out is that if one expresses  $W_{ij}(u)$  as a function w of  $Q_{0j}(u)$ ,

$$W_{\mu}(u) = W(Q_{0}(u))$$

then for long-tailed distributions,  $w(x) \sim \frac{1}{x}^2$ . By writing  $W_{\mu}$  as a function of  $Q_0$ , to an L-estimator one can form an equivalent iteratively reweighted estimator.

Given preliminary estimators  $\mu^*$  and  $\sigma^*$ , form  $\tilde{Q}_0^*(u) = \frac{\tilde{Q}(u) - \mu^*}{\sigma^*}$  and define  $\hat{\mu} = \frac{\int_0^1 w(\tilde{Q}_0^*(u)) \tilde{Q}(u) du}{\int_0^1 w(\tilde{Q}_0^*(u)) du}$ 

This estimator is a weighted average of Q with weights a function only of the size of the standardized residuals  $\tilde{Q}_0^*(u)$ .

For Student's t-distribution with m degrees of freedom,

$$J_0(u) = \psi (Q_0(u)) = \frac{m+1}{m} \frac{Q_0(u)}{1 + \frac{1}{m} Q_0^2(u)}$$

Consequently

w

$$W_{\mu}(u) = W_{\mu}(Q_0(u))$$
,  $W_{\sigma}(u) = W_{\sigma}(Q_0(u))$ 

with

$$\mu(\mathbf{x}) = \frac{m+1}{m} \frac{1 - (\mathbf{x}^2/m)}{\left[1 + (\mathbf{x}^2/m)\right]^2}, \quad \mathbf{w}_{\sigma}(\mathbf{x}) = \frac{m+1}{m} \frac{2\mathbf{x}}{\left[1 + (\mathbf{x}^2/m)\right]}$$

These windows deserve further investigation. However they appear to support the recommendation that robust estimators of location and scale may be obtained from preliminary estimators  $\mu^*$  and  $\sigma^*$  by the formulas (for a suitably chosen value of m )

$$\hat{\mu} = \frac{\int_{0}^{1} \tilde{Q}(u) \{1 + \frac{1}{m} \left(\frac{Q(u) - \mu^{*}}{\sigma^{*}}\right)^{2}\}^{-1} du}{\int_{0}^{1} \{1 + \frac{1}{m} \left(\frac{\tilde{Q}(u) - \mu^{*}}{\sigma^{*}}\right)^{2}\}^{-1} du}$$

$$\hat{\sigma}^{2} = \int_{0}^{1} \{\tilde{Q}(u) - \mu\}^{2} \{1 + \frac{1}{m} \left(\frac{\tilde{Q}(u) - \mu^{*}}{\sigma^{*}}\right)^{2} \}^{-1} du \left(\frac{m+1}{m}\right)$$

# 5. Variance and Influence Functions of Window Estimators

This section presents a conjectured formula for the asymptotic variance of a window estimator which is derived by representing it as an L-estimator

$$\hat{\mu} = \frac{\int_0^1 w(Q_0(u)) \ \tilde{Q}(u) \ du}{\int_0^1 w(Q_0(u)) \ du}$$

where  $w(x) = (1 + \frac{1}{m}x^2)^{-1}$  The question of deriving the theory of  $\mu$  as an M-estimator is open for research;  $\mu$  is an M-estimator if it satisfies

$$\int_0^1 \psi(\frac{\widehat{Q(u)} - \widehat{\mu}}{\sigma}) du = 0$$

for a suitable & function, here chosen to be

$$\Psi(\mathbf{x}) = \frac{\mathbf{x}}{1 + (\mathbf{x}^2/m)}$$

Under the assumption that the true quantile function is of the form  $Q(u) = \mu + \sigma Q_0(u)$ , and that  $Q_0(1 - u) = -Q_0(u)$ , signifying a symmetric distribution, we seek to find the variance V of the asymptotic distribution of  $\sqrt{n}(\mu - \mu)$ , which is normal with zero mean and asymptotic variance V.

From the asymptotic distribution theory of L-estimators

$$\mathbf{V} = \frac{\int_{0}^{1} |V(\mathbf{u})|^{2} d\mathbf{u}}{\{\int_{0}^{1} wQ_{0}(\mathbf{u}) d\mathbf{u}\}^{2}}$$

where

$$V'(u) = w(Q_0(u)) q(u) = \sigma w(Q_0(u)) q_0(u)$$

$$V(u) = \sigma v(Q_0(u)) ,$$

defining

۲

$$\mathbf{v}(\mathbf{x}) = \int^{\mathbf{x}} \mathbf{w}(\mathbf{y}) \, \mathrm{d}\mathbf{y} = \sqrt{m} \, \tan^{-1}(\mathbf{x}/\sqrt{m})$$

Further 
$$v(x)$$
 is the influence function of the estimator (Huber (1977), p. 17).

Note that for fixed x,  $v(x) \rightarrow x$  as  $m \rightarrow \infty$ . The formula for the variance of the robust estimator  $\mu$  may be written explicitly

$$\operatorname{Var}(\mu) = \frac{\sigma^2}{n} \quad \frac{\int_0^1 \left\{ \sqrt{m \tan^{-1} (Q_0(u)/\sqrt{m})} \right\}^2 du}{\left\{ \int_0^1 (1 + \frac{1}{m} Q_0^2(u))^{-1} du \right\}^2}$$

and can clearly be regarded as a generalization of the traditional formula for the variance of the sample mean. It is derived under the assumption of a symmetric but possibly long-tailed distribution.

To estimate  $Var(\mu)$  in practice, one might replace  $\sigma$  by  $\sigma$  and  $\tilde{Q}_0(u)$  by  $(\tilde{Q}(u) - \mu)/\sigma$  if a Quantile-Box plot of  $\tilde{Q}(u) - \mu$  indicates that it is symmetrically distributed about 0.

It should be noted that under the model  $Q(u) = \mu + \sigma Q_0(u)$ , with  $Q_0(1 - u) = -Q_0(u)$ ,  $\hat{\mu}$  estimates  $\int_0^1 w(Q_0(u)) Q(u) du \div \int_0^1 w(Q_0(u)) du = \mu$ while  $\hat{\sigma}^2$  estimates  $\sigma^2 \int_0^1 w(Q_0(u)) Q_0^2(u) du = \sigma^2 \int_0^1 Q_0^2(u) (1 + \frac{1}{m} Q_0^2(u))^2 du$ .

## 6. Non-parametric Probability Law Modeling

To interpret (as well as to form) location and scale parameters estimators from a data batch  $X_1, \ldots, X_n$  one must model its probability law. This section briefly mentions some new approaches which are currently being developed for non-parametric probability law modeling (see Parzen (1978)). They all involve both graphical and numerical analysis of the sample quantile function  $\tilde{Q}$  to find smoothing functions  $\tilde{Q}$ .

Quantile Box-Plots are introduced in the next section.

Quantile Residual Brownian Bridge Test. To say that the true quantile function Q(u) obeys the hypothesis  $H_0$ : Q(u) =  $\mu + \sigma Q_0(u)$  is to say that one can find values  $\mu$  and  $\sigma$  such that  $Q(u) = \mu + \sigma Q_0(u)$ fits  $\tilde{Q}$ . The fit of  $\tilde{Q}$  to  $\tilde{Q}$  can be judged by displaying the <u>quantile</u> residuals

$$R(u) = f_0 Q_0(u) \{ Q(u) - Q(u) \}, \quad 0 \le u \le 1$$

where  $f_0Q_0(u) = f_0(Q_0(u))$  is the <u>density-quantile</u> function corresponding to  $F_0$ . Under the null hypothesis ( $\sqrt{n}/\sigma$ ) R(u),  $0 \le u \le 1$  is asymptotically distributed as a stochastic process B(u),

 $0 \le u \le 1$  which is a modified Brownian Bridge process in the sense that its covariance kernel  $E(B(u_1)B(u_2))$  is not min  $(u_1, u_2) - u_1u_2$  but is modified due to the estimation of the parameters  $\mu$  and  $\sigma$ . To test whether the sample path R(u) looks like a sample path from a modified Brownian Bridge process one could use various functionals whose asymptotic distribution is known from their role in the conventional theory of Goodness of Fit Tests. The sample process traditionally considered for goodness of fit tests is

$$\widetilde{D}_{0}(u) = F_{0}((\widetilde{Q}(u) - \mu)/\sigma)$$

To estimate  $\sigma$  (needed in the asymptotic distribution of R(u)) one could use a non-parametric estimator such as

$$\widetilde{\sigma}_0 = \int_0^1 f_0 Q_0(u) d\widetilde{Q}(u) = \int_0^1 J_0(u) \widetilde{Q}(u) du$$

To estimate  $\mu$  and  $\sigma$  needed to form Q(u), one could use quick and dirty estimators  $\mu^*$  and  $\sigma^*$  formed from Quantile Box-Plots, or one could use asymptotically efficient estimators formed from <u>regression</u> analysis of the continuous process Q(u) (see section 8).

<u>Cumulative Weighted Spacings Brownian Bridge Tests</u>. To test whether the true quantile function Q(u) is of the form  $Q(u) = \mu + \sigma Q_0(u)$ , one need not first estimate  $\mu$  and  $\sigma$ . Instead, following Parzen (1978), form

$$\widetilde{D}(u) = \frac{1}{\sigma_0} \int_0^u f_0 Q_0(t) d\widetilde{Q}(t), \quad 0 \le u \le 1,$$

which is an estimator of

$$D(u) = \frac{1}{\sigma_0} \int_0^u f_0 Q_0(t) \, dQ(t) \quad 0 \le u \le 1 \quad \text{defining}$$
  
$$\sigma_0 = \int_0^1 f_0 Q_0(u) \, dQ(u) \quad . \quad \text{Under the null hypothesis,} \quad D(u) = u \quad , \text{ and it is}$$

conjectured that  $\sqrt{n} \{ D(u) - u \}$ ,  $0 \le u \le 1$  is asymptotically distributed as a Brownian Bridge Stochastic process.

By suitably choosing the null hypothesis  $H_0$  and the standard density-quantile function  $f_0Q_0$ , one can test the goodness of fit of any specified probability law (normal, exponential, Weibull, Cauchy, etc.) to the data.

Density-Quantile Function Autoregressive Estimation. Parzen (1978) discusses autoregressive estimators d(u) of

$$d(u) = D'(u) = \frac{1}{\sigma_0} - \frac{f_0 Q_0(u)}{f Q(u)}$$

which can be used to form estimators of fQ(u) .

The density quantile function fQ can be estimated also by forming autoregressive smoothers  $D_0(u)$  of

$$\widetilde{D}_{0}(u) = F_{0}((\widetilde{Q}(u) - \mu)/\sigma)$$

with density

$$\widetilde{d}_{0}(u) = f_{0}((\widetilde{Q}(u) - \mu)/\sigma)\widetilde{q}(u)/\sigma$$

The autoregressive density  $\hat{d}_0(u) = \hat{D}_0'(u)$  is an estimator of

$$d_0(u) = \frac{f_0((Q(u) - \mu))\sigma}{fQ(u)\sigma}$$

## 7. Quantile Box-Plots Diagnostic Measures

Given a data batch  $X_1, \ldots, X_n$ , a successful approach to "display" of the data has been the box plot introduced by Tukey (1977). Five values from a set of data are conventionally used: the extremes, the upper and lower H-values (H is an abbreviation for hinges or quartiles), and the M-value (median). The basic configuration of the box-plot display is a vertical box of arbitrary width and length equal to the distance HH (defined as upper H-value minus lower H-value and called the H-spread). A solid line (called the M-line) is marked within the box at a distance MH above the lower end of the box (MH equals M minus lower H). Dashed lines are extended from the lower and upper ends of the box a distance equal to the distance of the extremes from the hinges. If one wants to indicate a confidence interval for the median, one might add a line perpendicular to the M-line at its midpoint, and of length  $\pm$  HH/ $\sqrt{n}$ . The box-plot described should be called an H-Box Plot, because by replacing H-values by other types of values (called E-values and D-values) one can consider E-Box Plots and D-Box Plots.

The H-values are most conviently defined as  $\tilde{Q}(0.25)$  and  $\tilde{Q}(0.75)$ , the 1/4 percentiles. The E-values are the 1/8 percentiles  $\tilde{Q}(0.125)$  and  $\tilde{Q}(0.875)$ . The D-values are the 1/16 percentiles  $\tilde{Q}(0.0625)$  and  $\tilde{Q}(0.9375)$ . The mid-summaries of a data batch are

 $\tilde{\mu}(p) = \frac{1}{2} \{ \tilde{Q}(1-p) + \tilde{Q}(p) \}, \quad 0 \le p \le 0.5 .$ 

Of particular interest are

$$\tilde{\mu}_{M} = \tilde{\mu}(0.5)$$
,  $\tilde{\mu}_{H} = \tilde{\mu}(0.25)$ ,  $\tilde{\mu}_{E} = \tilde{\mu}(0.125)$ ,  
 $\tilde{\mu}_{D} = \tilde{\mu}(0.0625)$ .

When  $H_0$ :  $Q(u) = \mu + \sigma Q_0(u)$  holds, and  $Q_0(1 - u) \equiv -Q_0(u)$ ,  $\tilde{\mu}$  is an approximately unbiased estimator of  $\mu$ .

The average of the extreme-values of the sample will be denoted  $\tilde{\mu}(0)$ . The closeness of  $\tilde{\mu}(0)$  to the other  $\tilde{\mu}$  values may indicate whether the data batch has a short tailed symmetric distribution such as the uniform.

The mid-spreads of a data batch are

 $\tilde{S}(p) = \tilde{Q}(1-p) - \tilde{Q}(p)$ ,  $0 \le p \le 0.5$ .

Given a specified standardized quantile function  $Q_0$ , the <u>mid-scales</u> are defined by

$$\tilde{\sigma}(p) = \tilde{S}(p) \div S_{0}(p)$$

where  $S_0(p) = Q_0(1-p) - Q_0(p)$  is the mid-spread of  $Q_0$ . When  $H_0$  holds,  $\tilde{\sigma}(p)$  is an approximately unbiased estimator of  $\sigma$ . Of particular interest are

 $\widetilde{\sigma}_{_{_{_{_{_{}}}}}}=\widetilde{\sigma}\left(0.\,25\right)\;,\quad \widetilde{\sigma}_{_{_{_{_{}}}}}=\widetilde{\sigma}\left(0.\,125\right)\;,\quad \widetilde{\sigma}_{_{_{_{}}}}=\widetilde{\sigma}(0.\,0625)\;.$ 

Quick and dirty estimators of  $\mu$  and  $\sigma$  are given by

$$\widetilde{\mu} * = \frac{1}{4} \{ \widetilde{\mu}_{M} + \widetilde{\mu}_{H} + \widetilde{\mu}_{E} + \widetilde{\mu}_{D} \}$$

$$\tilde{\sigma} * = \frac{1}{5} \{ \tilde{\sigma}_{H} + 2\tilde{\sigma}_{E} + 2\tilde{\sigma}_{D} \}$$

Diagnostic tests for the validity of  $H_0$  are obtained by testing for the equality of the various  $\tilde{\mu}$  and  $\tilde{\sigma}$  values. More quantitative diagnostic measures could be defined as follows:

SK EW(p) = 
$$\{\tilde{\mu}_{M} - \tilde{\mu}(p)\} \div \tilde{S}(p)$$
  
TAIL(p) =  $\log \{\tilde{S}(p) \div \tilde{S}(0.25)\}$   
TAIL<sub>0</sub>(p) =  $\log \{S_{0}(p) \div S_{0}(0.25)\}$   
TAIL<sub>\$\phi(p)\$</sub> =  $\log \{\$^{-1}(p) \div \$^{-1}(0.25)\}$ 

When SKEW(p) is not significantly different from zero, we consider the data batch to have a symmetric distribution.

When the data passes a SKEW test for symmetry, it is checked for normality by comparing TAIL(p) with TAIL $\Phi(p)$ ; TAIL(p) significantly larger than TAIL $\Phi(p)$  indicates a long-tail distribution, and TAIL(p) significantly shorter than TAIL $\Phi(p)$  indicates either a shorttailed distribution (especially a uniform) or possibly a bimodal distribution. A seven-number summary of a data batch is provided by its M, H, E and D values, which suffice to compute mid-summaries, mid-scales, SKEW, and TAIL measures. To find re-expressions (transformations) of the data which make it more normal, one needs only the seven-number summary of the re-expressed data batch which are easily found as re-expressions of the seven-number summary of the original data batch.

In addition to the analytical measures of the data, one should form a graphical display of the quantile function  $\tilde{Q}(u)$  as a function on the unit interval  $0 \le u \le 1$ ; the H, E, and D boxes are drawn superimposed.

Quantile-Box Plots enable the investigator to detect "non-ideal" aspects of data batches by testing the data for normality by tests which determine the directions in which data fails to be normal, such as (1) longtailed distribution, (2) outliers, (3) bimodal distribution, (4) non-symmetric distribution.

To check for <u>symmetry</u>, inspect the shape of Q(u) within the boxes, as well as compare mid-summaries and examine the SKEW diagnostic measures.

When the data passes the test for symmetry the question of whether it has a <u>normal or long-tailed distribution</u> is decided using the TAIL diagnostic measures. Small TAIL values may indicate bimodal distributions. Data sets with outliers may also yield small TAIL values.

If the graph x = Q(u) has points with sharp rises ("infinite" slopes), then the probability density has a zero and will therefore have two (or more)

modes. If the points of sharp rise lie inside the H-Box we suspect the presence of several distinct populations generating the single data batch. If points of sharp rise lie outside the E-Box, we suspect outliers (values to be discarded for robust estimation).

A mode in the probability density function is indicated in the graph  $\mathbf{x} = \widetilde{\mathbf{Q}}(\mathbf{u})$  by a point of inflection (with "finite" slope). A horizontal segment in the graph is interpreted to mean a very large probability density there.

# Location and Scale Parameter Estimation as Regression Analysis of Sample Quantile Process

One can consider estimators, denoted  $\mu_{p,q}$  and  $\sigma_{p,q}$ , which use the sample quantile function Q(u),  $p \le u \le q$ ; this is equivalent to using a restricted set of order statistics  $X_{(np)}, \ldots, X_{(nq)}$  or a <u>trimmed</u> sample. A compact derivation of such formulas is given by Parzen (1978) who gives the representation

$$\begin{bmatrix} \hat{\mu} \\ \mu \\ \mu \\ \rho, q \end{bmatrix} = \begin{bmatrix} I_{\mu\mu} & I_{\mu\sigma} \\ I_{\mu\sigma} & I_{\sigma\sigma} \end{bmatrix}^{-1} \begin{bmatrix} T_{\mu, p, q} \\ T_{\sigma, p, q} \end{bmatrix}$$

where

$$T_{\mu,p,q} = \int_{p}^{q} W_{\mu}(u) \widetilde{Q}(u) du + \widetilde{Q}(p) W_{\mu L}(p) + \widetilde{Q}(q) W_{\mu R}(q)$$

$$T_{\sigma,p,q} = \int_{p}^{q} W_{\sigma}(u) \widetilde{Q}(u) du + \widetilde{Q}(p) W_{\sigma L}(p) + \widetilde{Q}(q) W_{\sigma R}(q)$$

$$I_{\mu\mu} = \int_{p}^{q} W_{\mu}(u) du + W_{\mu L}(p) + W_{\mu R}(q)$$

$$I_{\mu\sigma} = \int_{p}^{q} W_{\sigma}(u) du + W_{\sigma L}(p) + W_{\sigma R}(q)$$

$$I_{\mu\sigma} = \int_{p}^{q} W_{\sigma}(u) du + W_{\sigma L}(p) + W_{\sigma R}(q)$$

$$\int_{\sigma\sigma} = \int_{P} \int_{P} \int_{\sigma} \int$$

The weight functions are expressed in terms of the <u>density-quantile</u> function  $f_0Q_0(u) = f_0(Q_0(u))$  and the <u>score</u> function

$$J_0(u) = -(f_0Q_0)'(u) = \frac{-f_0'(F_0^{-1}(u))}{f_0(F_0^{-1}(u))} = \psi(Q_0(u)) .$$

$$\begin{split} \mathbf{W}_{\mu}(\mathbf{u}) &= \mathbf{J}_{0}^{t}(\mathbf{u}) \ \mathbf{f}_{0} \mathbf{\Omega}_{0}(\mathbf{u}) \\ \mathbf{W}_{\sigma}(\mathbf{u}) &= \mathbf{J}_{0}(\mathbf{u}) + \mathbf{\Omega}_{0}(\mathbf{u}) \ \mathbf{W}_{\mu}(\mathbf{u}) \\ \mathbf{W}_{\mu\mathbf{L}}(\mathbf{p}) &= \mathbf{f}_{0} \mathbf{\Omega}_{0}(\mathbf{p}) \left[ \frac{1}{\mathbf{p}} \ \mathbf{f}_{0} \mathbf{\Omega}_{0}(\mathbf{p}) + \mathbf{J}_{0}(\mathbf{p}) \right] \\ \mathbf{W}_{\mu\mathbf{R}}(\mathbf{q}) &= \mathbf{f}_{0} \mathbf{\Omega}_{0}(\mathbf{q}) \ \left[ \frac{1}{1 - \mathbf{q}} \ \mathbf{f}_{0} \mathbf{\Omega}_{0}(\mathbf{p}) - \mathbf{J}_{0}(\mathbf{p}) \right] \\ \mathbf{W}_{\sigma\mathbf{L}}(\mathbf{p}) &= \mathbf{\Omega}_{0}(\mathbf{p}) \ \mathbf{W}_{\mu\mathbf{L}}(\mathbf{p}) - \mathbf{f}_{0} \mathbf{\Omega}_{0}(\mathbf{p}) \\ \mathbf{W}_{\sigma\mathbf{R}}(\mathbf{q}) &= \mathbf{\Omega}_{0}(\mathbf{q}) \ \mathbf{W}_{\mu\mathbf{R}}(\mathbf{p}) + \mathbf{f}_{0} \mathbf{\Omega}_{0}(\mathbf{q}) \end{split}$$

For normally distributed data,

$$f_{0}(\mathbf{x}) = \phi(\mathbf{x}) = \frac{1}{\sqrt{2\pi}} e^{-(1/2)\mathbf{x}^{2}}, \quad F_{0}(\mathbf{x}) = \Phi(\mathbf{x}) = \int_{-\infty}^{\mathbf{x}} \phi(\mathbf{y}) \, d\mathbf{y} ,$$

$$f_{0}Q_{0}(\mathbf{u}) = \phi \Phi^{-1}(\mathbf{u}) = \frac{1}{\sqrt{2\pi}} \exp \left[-\frac{1}{2} \left| \Phi^{-1}(\mathbf{u}) \right|^{2} \right],$$

$$J_{0}(\mathbf{u}) = \Phi^{-1}(\mathbf{u}), \quad J_{0}'(\mathbf{u}) = \left\{ \phi \Phi^{-1}(\mathbf{u}) \right\}^{-1},$$

$$W_{\mu}(\mathbf{u}) = 1, \quad W_{\sigma}(\mathbf{u}) = 2 \Phi^{-1}(\mathbf{u}) .$$

$$W_{\mu}(\mathbf{p}) = \phi \Phi^{-1}(\mathbf{p}) \left\{ \frac{1}{p} \phi \Phi^{-1}(\mathbf{p}) + \Phi^{-1}(\mathbf{p}) \right\}$$

28

-

$$W_{\sigma R}(p) = \Phi^{-1}(p) W_{\mu L}(p) - \phi \Phi^{-1}(p)$$

When q = 1 - p,  $I_{\mu\sigma} = 0$ ,

$$I_{\mu\mu} = 1 - 2p + 2W_{\mu L}(p)$$
$$I_{\sigma\sigma} = 2 \int_{p}^{q} |\phi^{-1}(u)|^{2} du + 2\phi^{-1}(p) W_{\sigma L}(p)$$

The estimator

$$\hat{\mu}_{p,q} = \frac{\int_{p}^{1-p} \widetilde{Q}(u) \, du + W_{\mu L}(p) \{ \widetilde{Q}(p) + \widetilde{Q}(1-p) \}}{1 - 2p + 2W_{\mu L}(p)}$$

is similar to the Winsorized mean (with trimming proportion p ).

Robust Maximum Likelihood Estimation of Mean and Variance of a Normal Distribution. We may be willing to assume that our data is more normal than longtailed, but the shape of the true distributions is deviating slightly from the assumed normal model due to "wrong" values in the data set. We propose the following exploratory data analysis for robust estimation of  $\mu$  and  $\sigma$  from normal data with possible "outliers." We suggest the name "robust maximum likelihood estimators" for these estimators.

For selected values of p (at least p = 0.05, 0.25, and 0.45), (1) compute the estimators  $\mu_{p, 1-p}$  and  $\sigma_{p, 1-p}$ , and (2) plot the residuals.

$$\tilde{Q}(u) - \tilde{Q}(u) = \tilde{Q}(u) - \tilde{\mu}_{p,1-p} - \sigma_{p,1-p} \Phi^{-1}(u)$$

multiplied by  $\phi \Phi^{-1}(u)$ . Their values over the interval  $p \le u \le 1 - p$  can be used to test the hypothesis  $H_0$ . The residuals over the tail intervals  $u \le p$  and  $u \ge 1 - p$  can be used to test for the presence of "wrong values." One estimates  $\mu$  and  $\sigma$  by those estimators corresponding to the lowest value of p for which one finds no "wrong values" over the tail intervals.

# 9. Robust Regression

Formulating the estimation of location and scale parameters as a problem of <u>weighted regression</u> of the sample quantile function, with weights a function of  $Q_0(u)$ , leads to the amazing conclusion that asymptotically efficient estimators of  $\mu$  and  $\sigma$  are obtainable numerically by iterating ordinary regression calculations!

The iteratively reweighted estimators  $\mu$  and  $\sigma$  are also the solutions to the problem of estimating  $\mu$  and  $\sigma$  in the following weighted least squares linear regression problem:

$$X_j = \mu + \epsilon_j$$

where  $\varepsilon_i$  are independent normal with mean 0 and variance satisfying

$$\operatorname{Var}(\varepsilon_{j}) = \sigma^{2} \frac{1}{w_{j}}, w_{j} = w(\varepsilon_{j}^{*}), \varepsilon_{j}^{*} = \frac{X_{j} - \mu^{*}}{\sigma^{*}}$$
 (\*)

A regression model can be written

$$Y_j = \beta_1 X_{1j} + \ldots + \beta_k X_{kj} + \epsilon_j$$
,

where  $\{\varepsilon_j\}$  are independent random variables with quantile function known up to a parameter  $\sigma$ 

$$Q_{\sigma}(u) = \sigma Q_{0}(u)$$

To robustly estimate the coefficients  $\beta_1, \ldots, \beta_k$ ,  $\sigma$  assume first  $Q_0(u) = \Phi^{-1}(u)$ , corresponding to normality, and by ordinary least squares linear regression obtain preliminary estimators  $\beta_1^*, \ldots, \beta_k^*$ ; then form residuals

$$\boldsymbol{\varepsilon_{j}}^{*} = (\boldsymbol{Y_{j}} - \boldsymbol{\beta_{l}}^{*} \boldsymbol{X_{lj}} + \ldots + \boldsymbol{\beta_{k}}^{*} \boldsymbol{X_{kj}}) \div \boldsymbol{\sigma_{j}}^{*}$$

The next stage of estimators  $\hat{\beta}_1, \ldots, \hat{\beta}_k$ ,  $\hat{\sigma}^2$  are taken to be the least squares linear regression estimators under the assumption that  $\epsilon_j$  have variances defined by (\*). This process is iterated to yield robust estimators (compare Huber (1977), p. 38, Algorithm W).

The long-tailed character of the residuals  $\varepsilon_j^*$  should also be examined, using Quantile-Box plots.

One might consider non-parametric non-linear regression of Y on  $X_1, \ldots X_k$ . A density-quantile approach to non-linear non-

parametric regression has been described by Parzen (1977), and is currently being investigated by Prof. J. P. Carmichael. It has been applied to time series analysis by Prof. M. Pagano. It provides means of checking whether robust estimation of variances and correlations is provided by robust estimation of linear regression coefficients.

## 10. Do we need a new definition of Statistics?

Can statistics be made a subject that provides intellectually exciting pastimes (for the young and the mature), is regarded as relevant

by the creative scientist, and is appealing as a career to the mathematically talented?

An important step in achieving these desirable (and I believe attainable) goals is to alter the perception of the sample mean and the sample variance in elementary statistical instruction. Introductory Statistics is regarded by almost all college students (even by mathematically talented students) as a very dull subject. Perhaps one reason is that students enter the course knowing about a mean and a variance and leave the course knowing only about a mean and a variance. That statistics is in fact a live and vibrant discipline can be communicated to the student by emphasizing that there are many ways to estimate mean and variance, and more generally location and scale parameters. I believe that the discipline of statistics can be made more "glamorous" if intellectually sound and demonstratively useful concepts of statistical data analysis and robust statistical inference are incorporated in introductory statistical instruction. The perspective which this paper proposes for interpreting robust statistical inference is equivalent to a proposal for the definition of statistics:

"Statistics is arithmetic done by the method of Lebesgue integration."

I realize this definition sounds unbelievable and may never sell to the introductory student. But at least statisticians should understand to what extent it is true. Perhaps it provides a basis for a new sect of statisticians.

Can we all agree that a basic problem of statistics is an arithmetical one: find the average  $\overline{X}$  of a set of numbers  $X_1, \ldots, X_n$ ? Even grade school students (in the U.S.A.) nowadays know the answer:

$$\overline{\mathbf{X}} = \frac{1}{n} (\mathbf{X}_1 + \mathbf{X}_2 + ... + \mathbf{X}_n)$$

In words: list the numbers, add them up, and divide by n. What should be realized is that the foregoing algorithm is the method of Riemann integration.

The method of Lebesgue integration finds  $\overline{X}$  by first finding the distribution function F(x) of the data, defined by F(x) = fraction of  $X_1, \ldots, X_n \le x$ ,  $-\infty < x < \infty$ . Then  $\overline{X}$  is found as the mean of this distribution function, defined by the integral

$$\overline{\mathbf{X}} = \int_{-\infty}^{\infty} \mathbf{x} \, \mathrm{dF}(\mathbf{x})$$
.

To use an analogy to count a sack of coins, first arrange the coins in piles according to their value (pennies, nickels, dimes, quarters, and halfdollars), then count the number of coins in each pile, determine the value of each pile, and finally obtain  $\overline{X}$  as the sum of the values of the piles, divided by n. The role of statistics is to find more accurate estimators of the true mean by fitting a smooth distribution function  $\widehat{F}(x)$  to  $\widetilde{F}(x)$ .

Still more insight (and fidelity to the truth) is obtained by displaying the sample <u>quantile</u> function  $\tilde{Q}(u) = \tilde{F}^{-1}(u) = \inf \{x : \tilde{F}(x) \ge u\}$ , and fitting smooth quantile functions  $\tilde{Q}(u)$  to  $\tilde{Q}(u)$ . Then one computes the "sample average" by

$$\hat{\mu} = \int_0^1 \hat{Q}(u) \, du = \int_0^1 W_{\mu}(u) \, \hat{Q}(u) \, du \div \int_0^1 W_{\mu}(u) \, du \quad .$$

In words, the "average" of a sample is a weighted average of the numbers in the sample <u>arranged in increasing order</u>, with the weight of a number depending on its rank. This is the essence of robust statistical data analysis; all the rest is commentary.

### References

# Huber, P. [1977]. <u>Robust Statistical Procedures</u>. Regional Conference Series in Applied Mathematics <u>27</u>. SIAM: Philadelphia.

- McGill, R., Tukey, J. W., Larsen, W. A. [1978]. "Variations of Box Plots," American Statistician, 32, 12-16.
- Mosteller, F. and Tukey, J. W. [1977]. Data Analysis and Regression, Addison Wesley: Reading, Mass.
- Parzen, E. [1977]. "Nonparametric Statistical Data Science (A Unified Approach Based on Density Estimation and Testing for "White Noise")." Technical Report No. 47, Statistical Science Division, State University of New York at Buffalo.
- Parzen, E. [1978]. "Nonparametric Statistical Data Modeling," Journal of the American Statistical Association.
- Tukey, J. W. [1977]. <u>Exploratory Data Analysis</u>. Addison Wesley: Reading, Mass.

## Appendix

EXAMPLES OF QUANTILE-BOX PLOTS TIPPETT'S WARP BREAK DATA (Compare box plots in McGill, Tukey, Larsen [1978]).

Fossil data from yellow Limestone formation of northwestern Jamaica (from Chernoff, H. (1973), "The Use of Faces to Represent Points in k-Dimensional Space Graphically," Journal of the American Statistical Association, <u>68</u>, 361-368).

Variables 2 and 6 have zeroes in fQ (rises in Q). Variables 3 and 4 have proability masses (flat stetches in Q). Variables 1 and 5 are candidates for re-expression (logarithm for 1, square root for 5). Variables 2 and 5 suffice to classify the observations.





\_\_\_\_\_





and the second state of the



SECURITY CLASSIFICATION OF THIS PAGE (When Date Entered)	
REPORT DOCUMENTATION PAGE	BEFORE COMPLETING FORM
1. REPORT NUMBER 2. GOVT ACCESSION NO	. 3. RECIPIENT'S CATALOG NUMBER
2 achnical Report No. ARO-4	
4. TITLE (and Sublitie)	5. TYPE OF REPORT & PERIOD COVERE
A Density-Quantile Function Perspective on	Technical
Robust Estimation /	6. PERFORMING ORG. REPORT NUMBER
7. AUTHOR(a)	S. CONTRACT OR GRANT NUMBER(.)
Emanuel Parzen	DA AG29-76-G-0239
. PERFORMING ORGANIZATION NAME AND ADDRESS	10. PROGRAM ELEMENT, PROJECT, TASK
Statistical Science Division	AREA & WORK UNIT NUMBERS
State University of New York at Buffalo	
Amherst, New York 14226	
11. CONTROLLING OFFICE NAME AND ADDRESS	12. REPORT DATE March 1978
	13. NUMBER OF PAGES
14. MONITORING AGENCY NAME & ADDRESS(II dilierent from Controlling Office)	15. SECURITY CLASS. (of this report)
	Unclassified
	ISA DECLASSIFICATION DOWNGBADING
Approved for public release; distribution unli	mited.
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unli</li> <li>.</li> <li>17. DISTRIBUTION STATEMENT (of the abetract entered in Block 20, 11 different in</li> <li>NA</li> </ul>	mited.
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unlight of the about a constraint of the Army position, unless so desi</li> </ul>	mited. om Report) strued as an official gnated by other authorized
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unli </li> <li>17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different in NA </li> <li>18. SUPPLEMENTARY NOTES The findings in this report are not to be cons Department of the Army position, unless so desi documents.</li></ul>	mited. om Report) strued as an official gnated by other authorized
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unlight of the abstract entered in Block 20, if different in NA</li> <li>18. SUPPLEMENTARY NOTES The findings in this report are not to be const Department of the Army position, unless so desi documents. 19. KEY WORDS (Continue on reverse elde if necessary and identify by block number Robust Statistics Box Plots</li></ul>	mited. om Report) strued as an official gnated by other authorized
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlight of the public release; distribution unlight of the abstract entered in Block 20, if different in NA 18. SUPPLEMENTARY NOTES The findings in this report are not to be const documents. 19. KEY WORDS (Continue on feveres elde if necessary and identify by block number Robust Statistics Box Plots Exploratory Data Analysis Robust Regre</li></ul>	mited. om Report) strued as an official gnated by other authorized
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlight of the public release; distribution unlight of the eduted in Block 20, 11 different fr </li> <li>17. DISTRIBUTION STATEMENT (of the edutred in Block 20, 11 different fr</li> <li>NA</li> <li>18. SUPPLEMENTARY NOTES The findings in this report are not to be conse Department of the Army position, unless so deside documents. 19. KEY WORDS (Continue on reverse elde 11 necessary and identify by block number Robust Statistics Box Plots Exploratory Data Analysis Robust Regree Parameter Estimation</li></ul>	mited. om Report) strued as an official gnated by other authorized
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlight of the public release; distribution unlight of the abetract entered in Block 20, if different in NA 18. SUPPLEMENTARY NOTES The findings in this report are not to be const documents. 19. KEY WORDS (Continue on reverse elde if necessary and identify by block number Robust Statistics Box Plots Exploratory Data Analysis Robust Regree Parameter Estimation Quantile Functions</li></ul>	mited. om Report) strued as an official gnated by other authorized ) ssion
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlight of the public release; distribution unlight of the electric entered in Block 20, 11 different in NA 18. SUPPLEMENTARY NOTES The findings in this report are not to be const Department of the Army position, unless so desit documents. 19. KEY WORDS (Continue on reverse elde if necessary and identify by block number Robust Statistics Box Plots Exploratory Data Analysis Robust Regre Parameter Estimation Quantile Functions 20. ABSTRACT (Continue on reverse elde if necessary and identify by block number parameter Estimation</li></ul>	mited. om Report) strued as an official gnated by other authorized ") ssion
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unlight of the abstract entered in Block 20, if different in NA</li> <li>17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different in NA</li> <li>18. SUPPLEMENTARY NOTES The findings in this report are not to be const Department of the Army position, unless so desi documents. 19. KEY WORDS (Continue on reverse elde if necessary and identify by block number Robust Statistics Box Plots Exploratory Data Analysis Robust Regree Parameter Estimation Quantile Functions 20. ABSTRACT (Continue on reverse elde if necessary and identify by block number This paper provides an overview to a new new reserve.</li></ul>	mited. om Report) strued as an official gnated by other authorized by ssion proach to statistic
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unlight of the abetraci entered in Block 20, if different in NA</li> <li>17. DISTRIBUTION STATEMENT (of the abetraci entered in Block 20, if different in NA</li> <li>18. SUPPLEMENTARY NOTES The findings in this report are not to be const Department of the Army position, unless so desi documents. 19. KEY WORDS (Continue on reverse elde if necessary and identify by block number Robust Statistics Box Plots Exploratory Data Analysis Robust Regre Parameter Estimation Quantile Functions 20. ABSTRACT (Continue on reverse elde if necessary and identify by block number This paper provides an overview to a new ge data analysis and parameter estimation which contained on the sector of the sector</li></ul>	mited. om Report) estrued as an official gnated by other authorized official ssion neral approach to statistica uld be called the quantila
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unlight of the solution of the solution, unless so desi documents.</li> <li>19. KEY WORDS (Continue on reverse elde if necessary and identify by block number Robust Statistics Box Plots Exploratory Data Analysis Robust Regree Parameter Estimation Quantile Functions</li> <li>20. AustRACT (Continue on reverse elde if necessary and identify by block number This paper provides an overview to a new ge data analysis and parameter estimation which confunction approach. The aims of descriptive statistics</li> </ul>	mited. mited. om Report) strued as an official gnated by other authorized ssion neral approach to statistica uld be called the <u>quantile</u> stics (to graphically
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unlight of the source of</li></ul>	mited. om Report) or Report) or Report) or Report) strued as an official gnated by other authorized of ssion of neral approach to statistica uld be called the <u>quantile</u> <u>stics</u> (to graphically Quantile Box plots of the
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlight of the abetract entered in Block 20, il different in NA </li> <li>17. DISTRIBUTION STATEMENT (of the abetract entered in Block 20, il different in NA</li> <li>18. SUPPLEMENTARY NOTES The findings in this report are not to be conse Department of the Army position, unless so desi documents. 19. KEY WORDS (Continue on reverse elde if necessary and identify by block number Robust Statistics Box Plots Exploratory Data Analysis Robust Regree Parameter Estimation Quantile Functions 20. AUSTRACT (Continue on reverse elde if necessary and identify by block number This paper provides an overview to a new ge data analysis and parameter estimation which co function approach. The aims of descriptive statistics paper generation of the data) are obtained by sample quantile function. The aims of "proodness"</li></ul>	mited. om Report) om Report) strued as an official gnated by other authorized ssion neral approach to statistica uld be called the <u>quantile</u> <u>stics</u> (to graphically Quantile-Box plots of the of fit <sup>m</sup> are obtained by
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unlight of the public release; distribution unlight of the ebetrect entered in Block 20, il different in NA</li> <li>18. SUPPLEMENTARY NOTES</li> <li>19. KEY WORDS (Continue on reverse elde if netenenty and identify by block number Robust Statistics Box Plots</li> <li>Exploratory Data Analysis Robust Regree Parameter Estimation Quantile Functions</li> <li>20. Austract (Continue on reverse elde if necessary and identify by block number This paper provides an overview to a new ge data analysis and parameter estimation which confunction approach. The aims of "goodness"</li> </ul>	mited. mited. om Report) estrued as an official gnated by other authorized official proach to statistical and be called the <u>quantile</u> stics (to graphically Quantile-Box plots of the of fit <sup>m</sup> are obtained by
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unlight of the abstract entered in Block 20, if different in NA</li> <li>17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different in NA</li> <li>18. SUPPLEMENTARY NOTES The findings in this report are not to be considered by block number of the Army position, unless so deside documents. 19. KEY WORDS (Continue on reverse elde if necessary and identify by block number Robust Statistics Box Plots Exploratory Data Analysis Robust Regree Parameter Estimation Quantile Functions 20. ADSTRACT (Continue on reverse elde if necessary and identify by block number of the Statistics Box Plots Exploratory Data Analysis Robust Regree Parameter Estimation Quantile Functions 20. ADSTRACT (Continue on reverse elde if necessary and identify by block number of the spaper provides an overview to a new ge data analysis and parameter estimation which confunction approach. The aims of descriptive statis summarize and display the data) are obtained by sample quantile function. The aims of "goodness (continue on the second parameter by the state of the second parameter by the second parameter by the second parameter by the state of the second parameter by the second parameter by the state of the second parameter by the state of the second parameter by the second paramete</li></ul>	mited. mited. or Report) estrued as an official gnated by other authorized possion ssion neral approach to statistica uld be called the <u>quantile</u> stics (to graphically Quantile-Box plots of the of fit <sup>m</sup> are obtained by pro- inued on page 2)
<ul> <li>16. DISTRIBUTION STATEMENT (of this Report)</li> <li>Approved for public release; distribution unlight of the abstract entered in Block 20, il different in NA</li> <li>17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, il different in NA</li> <li>18. SUPPLEMENTARY NOTES The findings in this report are not to be considered by block number of the Army position, unless so deside documents. 19. KEY WORDS (Continue on reverse elde if necessary and identify by block number Robust Statistics Box Plots Exploratory Data Analysis Robust Regree Parameter Estimation Quantile Functions 20. ANSTRACT (Continue on reverse elde if necessary and identify by block number This paper provides an overview to a new ge data analysis and parameter estimation which confunction approach. The aims of descriptive statistics is unmarize and display the data) are obtained by sample quantile function. The aims of "goodness (continue on the second second</li></ul>	mited. mited. or Report) etrued as an official gnated by other authorized possion ssion neral approach to statistic uld be called the <u>quantile</u> stics (to graphically Quantile-Box plots of the of fit <sup>m</sup> are obtained by inued on page 2)

## Unclassified

LECURITY CLASSIFICATION OF THIS PAGE (When Date Entered)

## page 2

fitting smooth quantile functions to the sample quantile function. The aims of <u>parameter estimation</u>, especially robust estimation of location and scale parameters, are attained by regression analysis of the sample quantile function. (The goal of a statistician in analyzing a batch of data  $X_1^{\prime\prime}, \ldots, X_n^{\prime\prime}$  should be both "estimation of parameters" and "goodness of fit." By "goodness of fit" is meant fitting of the observed sample probabilities by

a smooth probability law.)

Quantile functions are defined in Section 2. Window estimators of location and scale parameter estimation are defined in Section 3; their equivalence to L-estimators is discussed in Section 4. A conjectured expression is given in Section 5 for the asymptotic variance of window estimators. New approaches being developed for non-parametric probability law modeling are mentioned in Section 6; quantile box-plots are introduced in Section 7. Section 8 discusses location and scale parameter estimation using trimmed samples. Robust regression is the subject of Section 9. A new definition of statistics is proposed in Section 10.

and the second sec

#### Unclassified

SECURITY CLASSIFICATION OF THIS PAGE(When Dete Entered)