

(i)

5 Oct 2001

Resistant

In insensitive to localized misbehavior in the data

eg. median

~~mean~~

Robust

In insensitive to departures from assumptions surrounding an underlying probability model

mean

trimmed mean

ii

5 Oct 2001

M-estimates

The Biweight

$$y^* = \frac{\sum w_i y_i}{\sum w_i}$$

$$w_i = \begin{cases} \left(1 - \left(\frac{y_i - y^*}{cS}\right)^2\right)^2 & \left(\frac{y_i - y^*}{cS}\right)^2 < 1 \\ 0 & \text{otherwise} \end{cases}$$

$$S = \text{median}\{|y_i - y^*|\}$$

$$\approx \frac{1}{2} \text{IQR}$$

S estimates " $\frac{2}{3}\sigma$ "

eg $c=6$, residuals count up to 4σ

Iterate to convergence

①

12 Oct. 21

The heuristic of robust regression. Consider a weight function like

$$w(r) = (1-r^2)^2 \quad |r| \leq 1$$

$$= 0 \quad \text{otherwise}$$

Consider a preliminary estimate $\hat{\beta}^-$ and residuals $\hat{\epsilon}_i^- = y_i - X_i \hat{\beta}^-$

Regress y_i on X_i with weight $w\left(\frac{\hat{\epsilon}_i^-}{\sigma}\right) = w_i$

ie. $\hat{\beta}$ to satisfy

$$\sum_i (y_i - X_i \hat{\beta}) w_i X_i = 0 \quad \text{cp.}$$

In the limit

$$\sum_i (y_i - X_i \hat{\beta}) w\left(\frac{y_i - X_i \hat{\beta}}{\sigma}\right) X_i = 0 \quad \text{cp. } (y - X\hat{\beta})^T w X = 0$$

ie.

$$\psi(r) = r w(r) \quad \text{or} \quad \boxed{w(r) = \psi(r)/r}$$

$$= r(1-r^2)^2 \quad |r| \leq 1$$

$$= 0 \quad \text{otherwise}$$

$$\rho(r) = \int_0^r \frac{1}{2}(1-s^2)^2 ds^2 = -\int_0^{r^2} \frac{1}{2}(1-t)^2 d(1-t) \quad t = s^2$$

$$= \int_{r^2}^1 \frac{1}{2} u^2 du = \frac{1}{2} \left. \frac{u^3}{3} \right|_{r^2}^1 = \frac{1}{6} [1 - (1-r^2)^3] \quad \frac{1}{2}(1-r^2)^2 2r$$

$$\boxed{\rho(r) = \frac{1}{6} [1 - (1-r^2)^3]} \quad \begin{matrix} |r| \leq 1 \\ = \frac{1}{6} \quad |r| > 1 \end{matrix}$$

(2)

7 Oct. 01

This approach suggests the approximate variance matrix

$$\text{var}(\hat{\beta} - \beta) \sim \sigma^2 (X^T W X)^{-1}$$

e.g. $\hat{\sigma} = \frac{3}{4} \text{IQR}$, median $\{|\hat{\epsilon}_i|\} = 6745$

Asymptotic normality from M-estimator theory

Start eg. with L_1 estimator

3

2 Oct. 01

Robust / Resistant Estimates

Would like an automatic way to handle outliers and long-tailed distributions.

Give up full efficiency, eg. in normal case, but obtain protection against outliers and nonnormality.

$$y_i = \underset{\sim}{x}_i^T \underset{\sim}{\beta} + \epsilon_i$$

M-estimate

$$\min_{\underset{\sim}{\beta}} \sum_{i=1}^n \rho \left(\frac{y_i - \underset{\sim}{x}_i^T \underset{\sim}{\beta}}{\sigma} \right)$$

or to satisfy

but need

$$\sum_{i=1}^n \psi \left(\frac{y_i - \underset{\sim}{x}_i^T \underset{\sim}{\beta}}{\sigma} \right) \underset{\sim}{x}_i = \underset{\sim}{0}$$

$$\psi = \rho'$$

(4)

Chapter 8 in $\mathbb{K} \times \mathbb{R}$

How to compare?

→ Oct 01

Examples of ρ :

1. OLS

$$\rho(r) = r^2$$

2. L_1

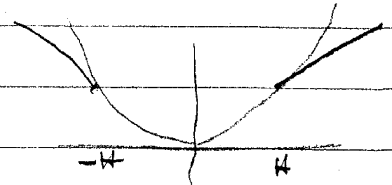
$$\rho(r) = |r|$$

3. Huber

$$\rho(r) = r^2$$

$$|r| \leq H$$

$$= H(2|r| - H)$$



H : large: normal

H small: L_1

$H = 1.345$ gives 95% efficiency at normal
 ρ : convex so math convenient

4. Bisquare

$$\rho(r) = \frac{B^2}{6} \left[1 - \left(1 - \left(\frac{r}{B} \right)^2 \right)^3 \right]$$

$$|r| \leq B$$

$$= \frac{B^2}{6} \text{ otherwise}$$

(2)

Example. Detecting circular arcs in image. Roads. 8 Oct. 2001

Fitting a circular arc when outliers severely perturb fit

$$C: (x-a)^2 + (y-b)^2 = R^2$$

$$x^2 + y^2 = (-2a)x + (-2b)y + (a^2 + b^2 - R^2)$$

$$z = \theta_1 x + \theta_2 y + \theta_3$$

CAE: circular arc estimator

DRO: Duda road operator
- nonlinear line detector

Theil-Ser: median of $\binom{N}{2}$ slopes

RM: repeated median

LSM: least median of squared residuals

MM: variant of M-estimator
- estimates scale too

SW: sliding window

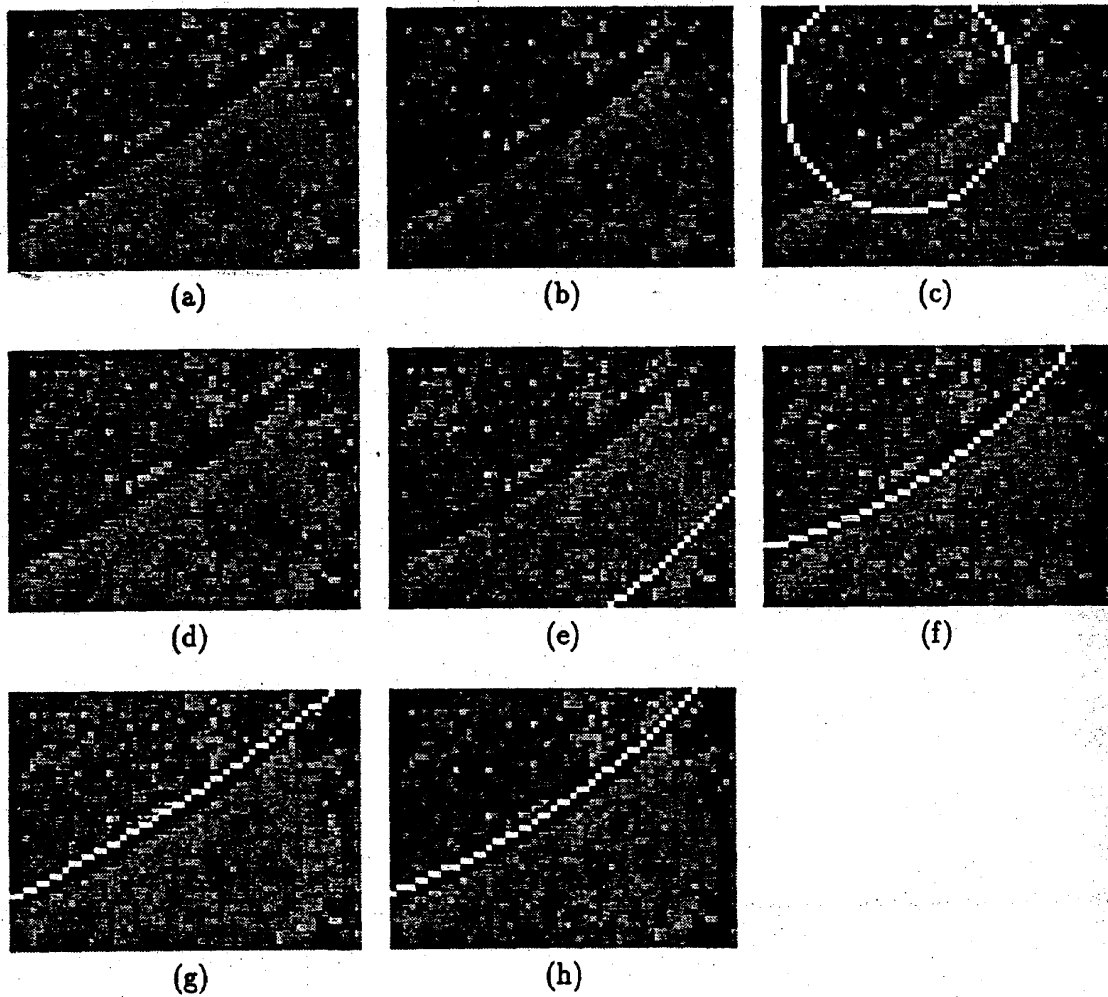


Figure 15: Data set 2: (a) Input image; (b) DRO output; (c) OLS fit; (d) Theil-Sen CAE fit (not contained in the displayed window); (e) RM CAE fit; (f) LMS fit; (g) nonlinear MM fit; (h) SW fit using an unbiased OLS estimator.

18 Oct. 01

Robust estimates.

M-estimates. Return to the regression problem

$$y \sim X\beta + \varepsilon$$

Consider the problem

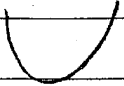
$$\min_{\beta} \sum_{j=1}^n \rho\left(\frac{y_j - x_j^T \beta}{s}\right) \quad x_j \text{ } j\text{-th row}$$

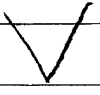
s : a known or estimated scale parameter, eg,

$$1.48 \left(\text{med}_j \left| y_j - x_j^T \hat{\beta}_{OLS} - \text{med}_j \left(y_j - x_j^T \hat{\beta}_{OLS} \right) \right| \right)$$

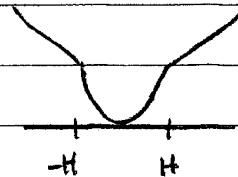
1.48 works for normal

Here $\rho(r)$ is a (robust) loss function

eg. $\rho(r) = r^2$ OLS 

$\rho(r) = |r|^p$ L_p regression  $p=1$

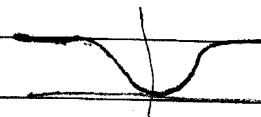
$\rho(r) = r^2 \quad |r| \leq H$
 $= H|r| - \frac{1}{2}H^2 \quad |r| > H$ Huber



(2)

8 Oct 01

$$\rho(r) = \begin{cases} \frac{B^2}{6} \left[1 - \left(1 - \frac{r}{B} \right)^2 \right]^3 & |r| \leq B \\ \frac{B^2}{6} & |r| > B \end{cases}$$



biweight
bisquare

Might pick tuning constants H, B, \dots so
95% efficient in the normal error case.

Note, 1. The location problem is a particular
case with $\tilde{X} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

L_2 regression gives sample mean
 L_1 " " " " median

2. \exists other rules for estimating scale

Properties of $\rho(\cdot)$

- i) $\rho(r) \geq 0$, $\rho(r)$ non decreasing for $r \geq 0$
- ii) $\rho(0) = 0$
- iii) $\rho(-r) = \rho(r)$
- iv) $\rho(\cdot)$ continuous for all but a finite number of r

(3)

8. Oct. 01

Differentiate wrt β and write $\psi = \rho'$

$$\sum_{j=1}^n x_j^T \psi\left(\frac{y_j - x_j^T \hat{\beta}}{\sigma}\right) = 0 \quad (*)$$

A set of nonlinear equations

- solve iteratively
 - good starting values needed
- L₁ (linear programming)

Write $w(r) = \psi(r)/r$

(*) becomes

$$\sum_{j=1}^n w\left(\frac{y_j - x_j^T \hat{\beta}}{\sigma}\right) x_j^T (y_j - x_j^T \hat{\beta}) = 0$$

OR

$$\sum_{j=1}^n w_j x_j^T (y_j - x_j^T \hat{\beta}) = 0$$

with data dependent weights

$$w_j = w\left(\frac{y_j - x_j^T \hat{\beta}}{\sigma}\right)$$

IRLS regress $\sqrt{w_j} y_j$ on $\sqrt{w_j} x_{ij}$

until "convergence"

w(r)

2

1/|r|

1

H/|r|

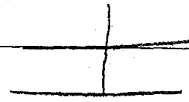
$(1 - (r/B)^2)^2$

0

$|r| \leq H$
 $|r| > H$

$|r| \leq B$

$|r| > B$



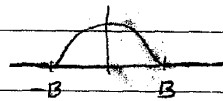
OLS



L1



Huber



Bisquare

Tri-square

Notes

- 1. ρ : nondecreasing \rightarrow estimate exists and unique
- 2. $\rho(r) = \text{constant}$ for $|r| > \text{constant}$
 \Rightarrow possible deletion of observations

Classification

$\psi(r) = 0$ for $|r|$ sufficiently large, eg. Bisquare
hard redescender

$\psi(r) \rightarrow 0$ as $|r| \rightarrow \infty$, eg. Cauchy
soft redescender $w(r) = 1 / (1 + (r/c)^2)$

Huber

nondecreasing

(5)

8 Oct 01

Biweight estimate of location

$$\hat{\mu} = \frac{\sum_i w_i y_i}{\sum_i w_i}$$

$$w_i = \left(1 - \left(\frac{y_i - \hat{\mu}}{c s}\right)^2\right)_+$$

$$s = \text{med}\{|y_i - \hat{\mu}|\}$$

$$c = 6, 9$$

Seems to converge quickly

Freedman and Diaconis (1982) Ann Statist 10, 454-461
inconsistency results

(6)

8 Oct 01

Heuristics in location case (Jeffries)

y_i from $f\left(\frac{y-\mu}{\sigma}\right) \frac{1}{\sigma}$ σ : "known"

log likelihood $\sum_i \log f\left(\frac{y_i-\hat{\mu}}{\sigma}\right)$

$$\frac{\partial}{\partial \mu}: -\frac{1}{\sigma} \sum_i \frac{1}{f\left(\frac{y_i-\hat{\mu}}{\sigma}\right)} f'\left(\frac{y_i-\hat{\mu}}{\sigma}\right) = 0$$

$$\text{or} \quad \sum_i w_i (y_i - \hat{\mu}) = 0$$

$$w_i = \frac{1}{f\left(\frac{y_i-\hat{\mu}}{\sigma}\right)} \cdot \frac{f'\left(\frac{y_i-\hat{\mu}}{\sigma}\right)}{\left(y_i-\hat{\mu}\right)}$$

$$\text{Cauchy: } f = \frac{1}{1+y^2}, \quad f' = -\frac{2y}{(1+y^2)^2}$$

$$w = -\frac{1}{f} \frac{f'}{y} = \frac{(1+y^2) \cdot 2y}{(1+y^2)^2} \cdot \frac{1}{y} = \frac{2}{(1+y^2)}$$

Jeffries took mixture of normal and uniform

⑦

8 Oct. 01

Approximate distributionSuppose $\underline{y} = \underline{X} \underline{\beta} + \underline{\varepsilon}$ distribution $\frac{\varepsilon_j}{\sigma} : p$

$$\hat{\underline{\beta}} \sim N_p \left(\underline{\beta}, \frac{E_p(\psi^2)}{[E_p(\psi')]^2} (\underline{X}^T \underline{X})^{-1} \sigma^2 \right)$$

Estimate of the covariance matrix

$$\frac{(ms)^2}{n-p} \frac{\sum_{j=1}^n \psi \left(\frac{\hat{\varepsilon}_j}{s} \right)^2}{\left[\sum_{j=1}^n \psi' \left(\frac{\hat{\varepsilon}_j}{s} \right) \right]^2} (\underline{X}^T \underline{X})^{-1} \quad \psi(r) = r w(r)$$

Sometimes $s^2 (\underline{X}^T \underline{W} \underline{X})^{-1}$

$$\text{e.g. } s^2 = \frac{3}{4} \text{IQR}$$

8

8 Oct 01

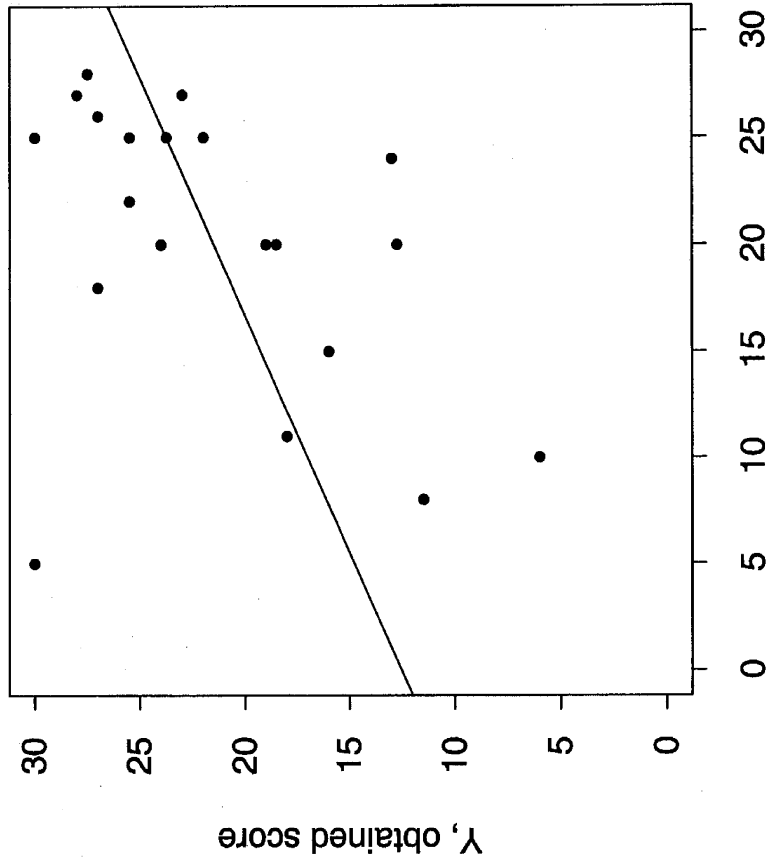
In practice run OLS and rcb/res in parallel.

If results dissimilar, indication on an invalid assumption (of classical regression)

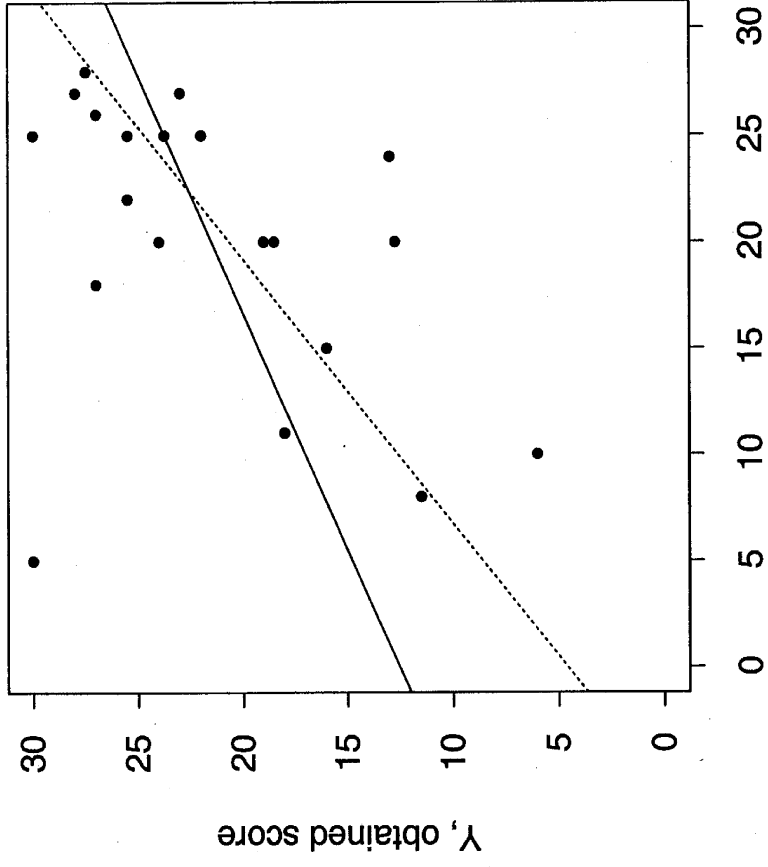
Look at residuals for possible outliers
eg. index plot

Look at index plot of weights

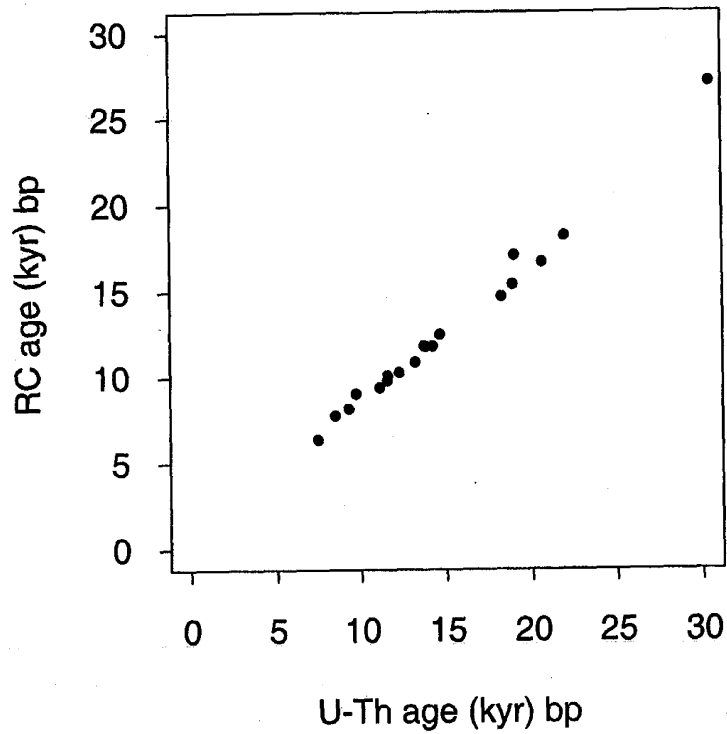
Stat 131a: Midterm



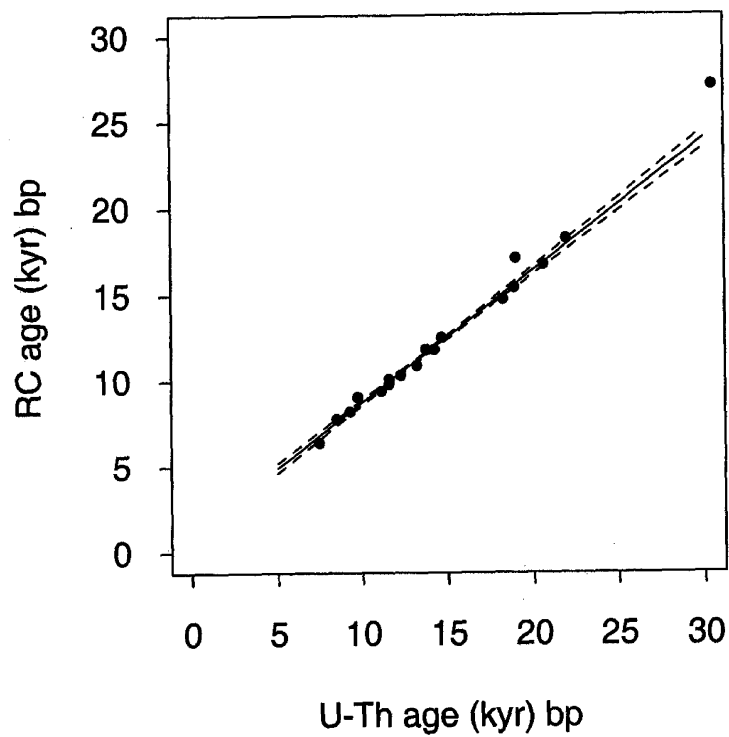
Stat 131a: Midterm + outlier



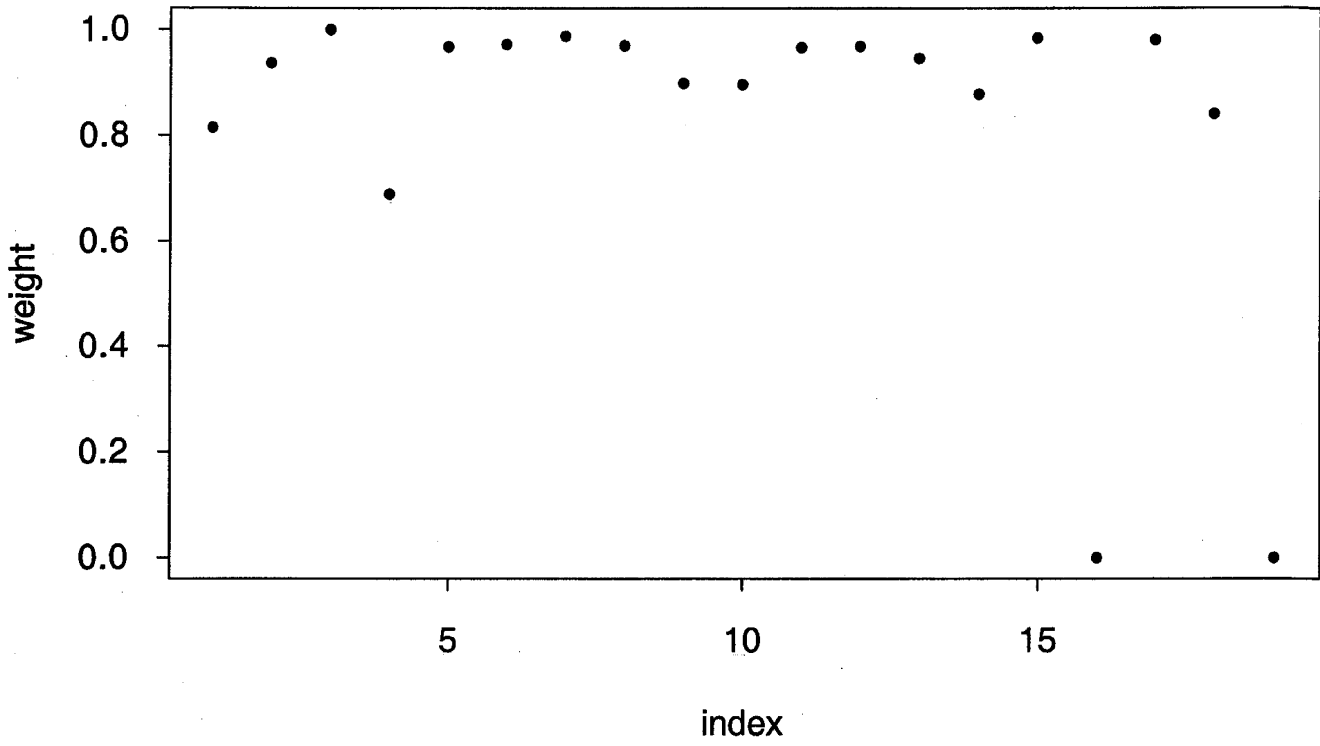
Radiocarbon age vs. uranium-thorium age



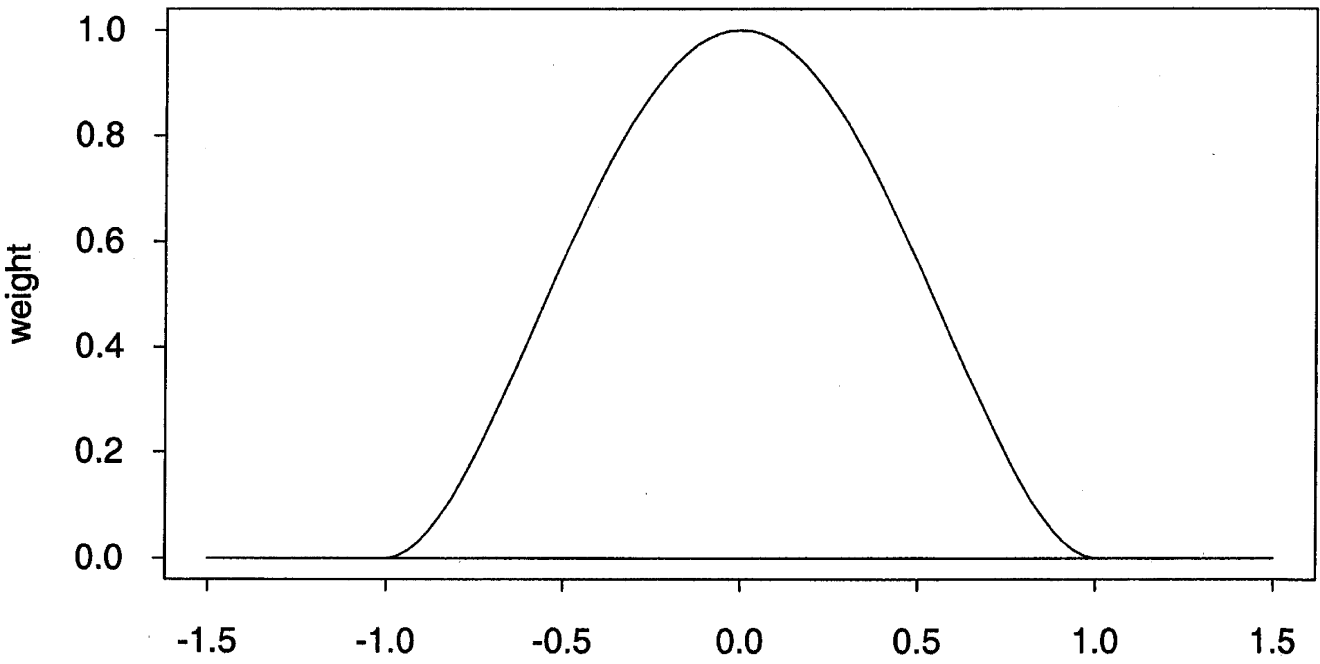
Radiocarbon age vs. uranium-thorium age



Weights plot



Bisquare weights



Summary of Robust/Resistant Methods

1. Can lead to identification of outliers
2. Can suggest where present model does not apply.
3. Leads to parameter estimates which are not sensitive to arbitrary changes in any small part of the data
 - breakdown point

Discuss anomalies with subject matter experts

10

8 Oct. 01

Functions in Splus

location.m ()

bisquare, huber

l1fit ()

L_1

rreg ()

converged Huber followed by
bisquare

rbiwt ()

bisquare, single α

ltsreg ()

trunc half $|\hat{\epsilon}_i|$

robust ()

glm (, family = robust)

gam ()