

Bradley Efron

# STATISTICAL DATA ANALYSIS

IN the Computer Age

B. Efron & R. Tibshirani

"Science" July 26 1991  
p390-395

"Today's Data analyst can afford to expend more computation on a single problem than the world's yearly total of statistical computation in the 1920's."

$n = 15$  Observed Lifetimes:

$.143$     $.182$     $.256$     $.260$     $.270$     $.437$   
•—————•—————•—————•—————•—————•  
 $x_1$     $x_2$     $x_3$     $x_4$     $x_5$     $x_6$

$.509$     $.611$     $.712$     $1.04$     $1.09$     $1.15$   
•—————•—————•—————•—————•—————•  
 $x_7$     $x_8$     $x_9$     $x_{10}$     $x_{11}$     $x_{12}$

$1.46$     $1.88$     $2.08$   
•—————•—————•  
 $x_{13}$     $x_{14}$     $x_{15}$

Sample Mean  
 $\bar{x} = .804$

Estimated Standard Error  
 $\hat{\sigma} = .155$

$$\sqrt{\left[ \frac{\sum (x_i - \bar{x})^2}{n-1} \right]^{\frac{1}{2}}}$$

Sample Median  
 $\hat{\theta} = .611$

$\hat{\sigma} = ???$

We have

DATA  $y$   
 $\approx$

e.g.  $y = (x_1, x_2, \dots, x_{15})$  random  
sample  $\approx$  of 15 lifetimes

PARAMETER OF Interest  $\theta$

e.g. the true expected lifetime  
or true median lifetime.

We Want to estimate

$\theta$  from  $y$   
 $\approx$

# Two Basic Questions

Question 1: What statistic  $\hat{\theta}(y)$  should we use to estimate  $\theta$ ?

Question 2: How accurate is  $\hat{\theta}$  as an estimate of  $\theta$ ?

# Maximum Likelihood Theory

Answer 1: Use  $\hat{\theta}(y)$   
the MLE

Answer 2: Standard Error  
of  $\hat{\theta}$  is approximately

$$\hat{\sigma} = \frac{1}{\sqrt{\text{Fisher Info}}}$$

BOOTSTRAP is a more general way to answer Q2.

- Less Parametric Modelling  
(even nonparametric)

- More Computation  
(x100 or 1000)

- Automatic  
(Algorithm)

# The Simplest Situation

$F$   $\xrightarrow{\text{Random Sample}}$   $(x_1, x_2, \dots, x_n) = \underline{y}$   
true distribution

Statistic of Interest:  $\hat{\theta}(\underline{y}) \rightarrow \bar{X} = \sum_{i=1}^n x_i / n$

Simple formula for standard error:  
 $\sigma(F) = [\mu_2(F) / n]^{1/2} \quad \star$

Where

$$\begin{aligned} \mu_2(F) &= 2^{\text{nd}} \text{ Central Moment of } F \\ &= E_F [X - E_F \{X\}]^2 \\ &= \text{Variance of one observation "X"} \end{aligned}$$



# Estimating $\sigma(F)$

EMPIRICAL  
Distribution:  $\hat{F}$ :  $\frac{1/n}{x_1} \quad \frac{1/n}{x_2} \quad \frac{1/n}{\dots} \quad \dots \quad \frac{1/n}{\dots} \quad \frac{1/n}{x_n}$

$$\mu_2(\hat{F}) = \sum_{i=1}^n (x_i - \bar{x})^2 / n$$

Gives estimated standard error

$$\begin{aligned} \sigma(\hat{F}) &= [\mu_2(\hat{F}) / n]^{1/2} \quad \star \\ \hat{F} \text{ for } F &= \left[ \sum_{i=1}^n (x_i - \bar{x})^2 / n^2 \right]^{1/2} \end{aligned}$$

a nice simple formula.

## For More Complicated Statistics

e.g.  $\hat{\theta}(y) = \text{Sample Median}$

- No Simple Formula for  $\sigma(F) = \text{Standard Error of } \hat{\theta}(y)$
- Can't "Substitute  $\hat{F}$  for  $F$ "
- Bootstrap: computer algorithm for finding numerical value of

$$\sigma(\hat{F}) \equiv \hat{\sigma}$$

Empirical  $\rightarrow$

$\leftarrow$  Bootstrap Estimate of Standard Error

# BOOTSTRAP SAMPLING

$$\text{" } \hat{F} \rightarrow \underset{\sim}{y}^* = (x_1^*, x_2^*, \dots, x_n^*) \text{"}$$

Means that you

"Draw a random sample  
of size  $n$  from  $\hat{F}$ "

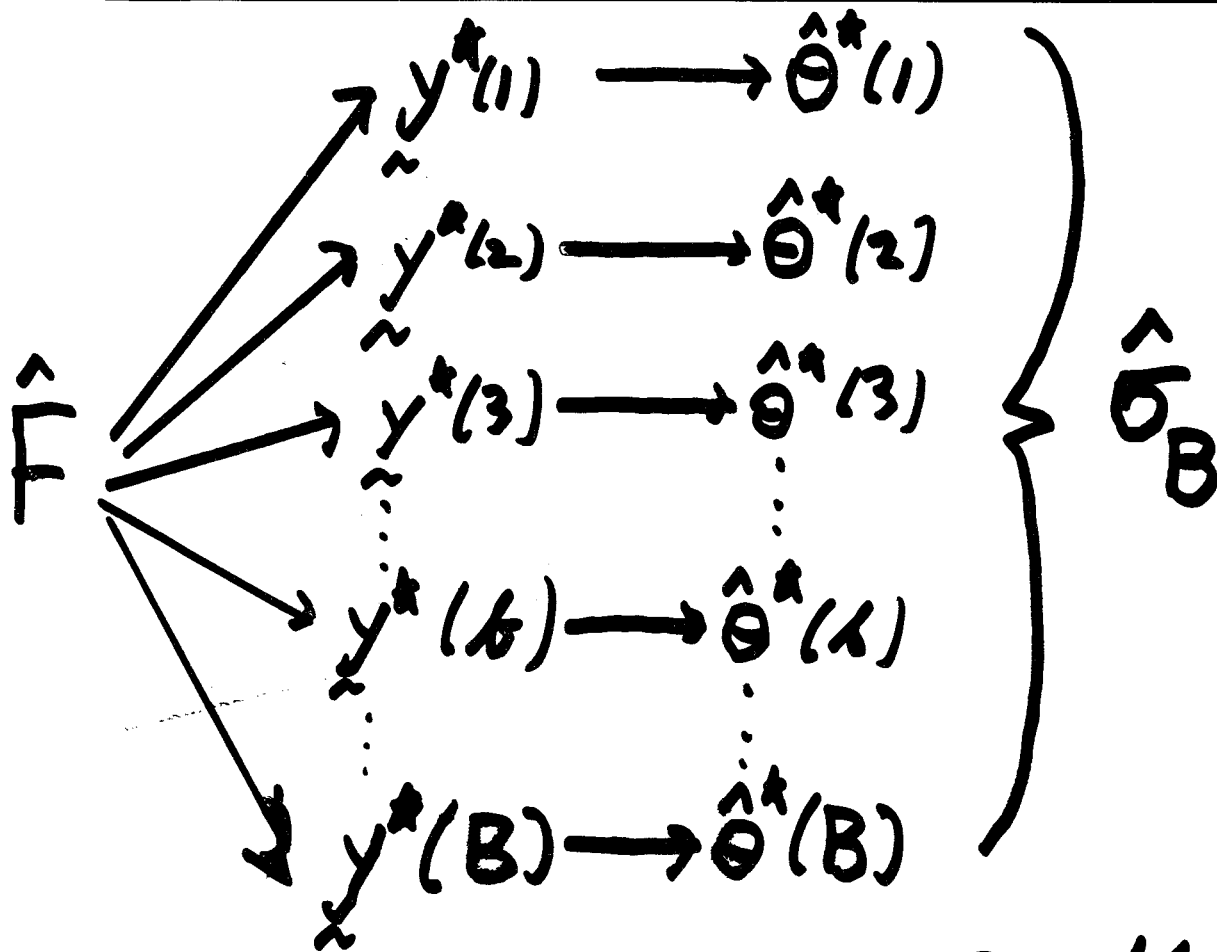
Equivalent:  $\underset{\sim}{y}^* = (x_1^*, x_2^*, \dots, x_n^*)$   
is a random sample of size  $n$   
drawn with replacement from

$$\{x_1, x_2, \dots, x_n\}$$

Definition:  $\underset{\sim}{y}^*$  is a Bootstrap Sample

# BOOTSTRAP ALGORITHM

---



$$\hat{\sigma}_B = \left\{ \frac{\sum [\hat{\theta}^*(i) - \hat{\theta}^*(1)]^2}{B-1} \right\}^{1/2}$$

As  $B \rightarrow \infty$ ,  $\hat{\sigma}_B \rightarrow \hat{\sigma}_0 = \sigma(\hat{F})$ ,  
 the bootstrap estimate of standard error.

# Results for the 15 lifetimes

B = 100 Bootstraps

Sample Mean  $\hat{\sigma}_{100} = .156$   $\hat{\sigma}_{20} = .155$   
 $\left[ \left\{ \sum (x_i - \bar{x})^2 / n \right\}^{1/2}$

Sample Median  $\hat{\sigma}_{100} = .229$

"boot (y, 100, "median")"

Usually  $B=100$  is plenty:

Lifetime Data,  $\hat{\theta} = \text{Sample Median}$

$B:$  25    50    100    200    1000

$\hat{\sigma}_B:$  .23    .22    .23    .25    .25

\* For Standard Errors

# MORAL

- Most of the familiar methods of getting standard errors are approx. versions of  $\hat{\sigma} = \sigma(\hat{F})$ .
- Bootstrap evaluates  $\sigma(\hat{F})$  directly by brute force Monte Carlo.
- Mathematical complications of step  $y \rightarrow \hat{\theta}$  don't bother the bootstrap.

# Ellipticity Data

• 51 Pairs  $x_i = (y_i, z_i)$   $i = 1, 2, \dots, 51$

•  $y_i = \text{temperature}$   $z_i = \text{ellipticity}$

• Model:

$$z_0 = \frac{a + b \cdot e^u}{1 + e^u}$$

where

$$u = d \cdot \left[ \frac{1}{c} - \frac{1}{y + 273.15} \right]$$

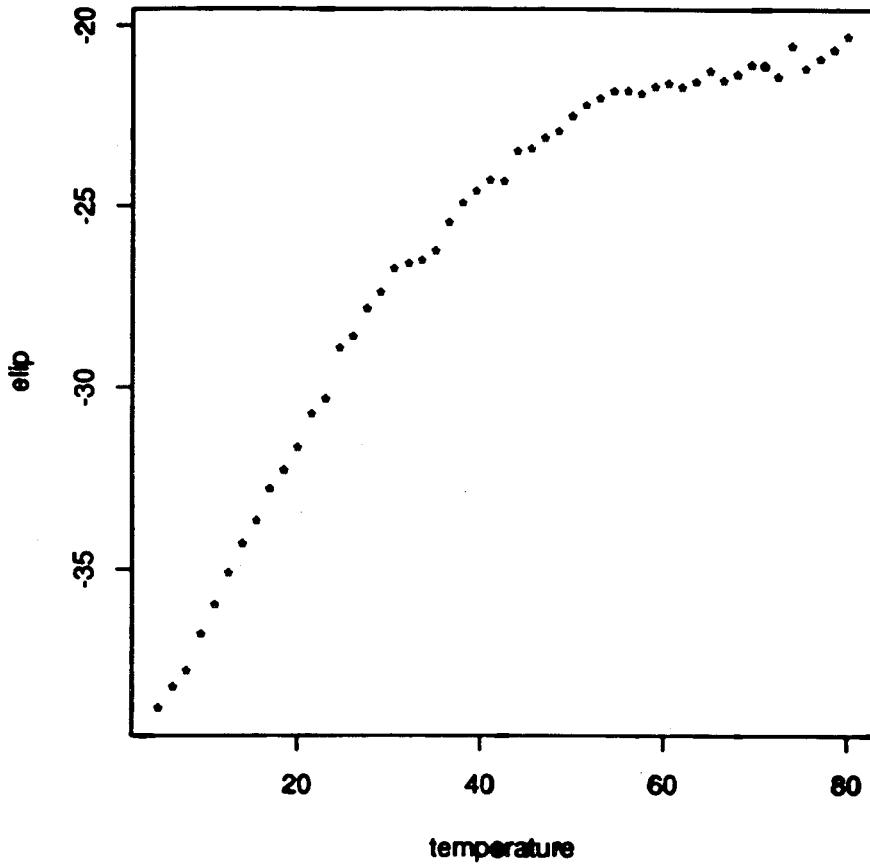
• Nonlinear model with four parameters  $(a, b, c, d)$

• Fit by least squares:

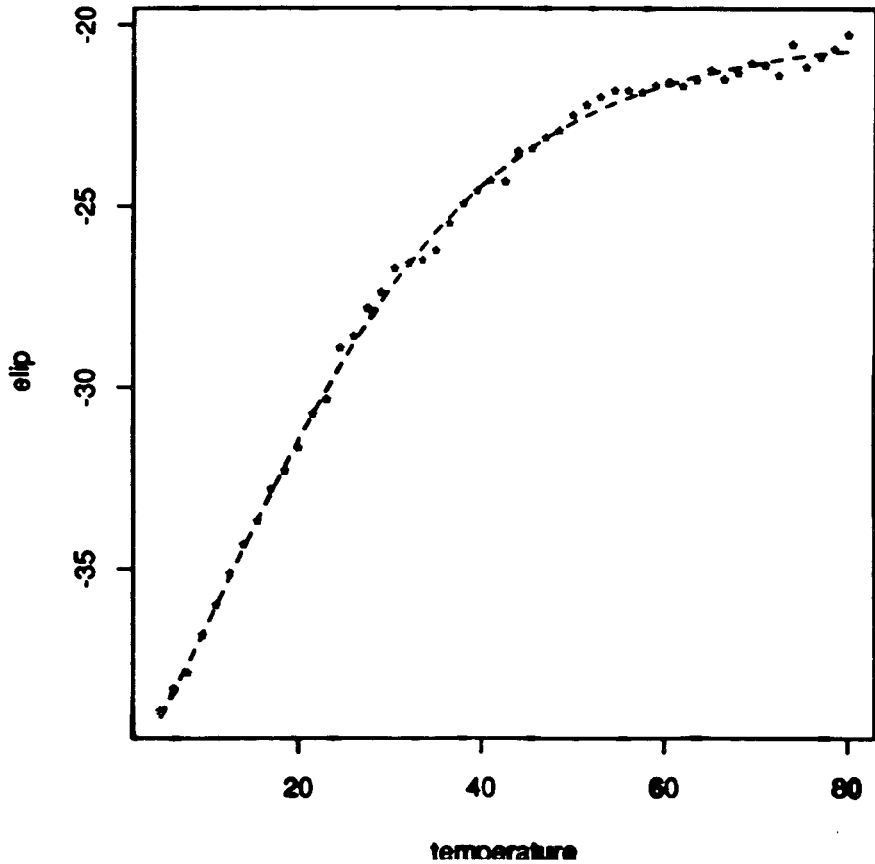
$$(a, b, c, d) \min \sum_{i=1}^{51} [z_i - z_{0i}]^2$$



51 pairs (temp, ellipticity)



51 pairs (temp, ellipticity)



$$\hat{a} = -48.3$$

$$\hat{b} = -20.1$$

$$\hat{c} = 287.5$$

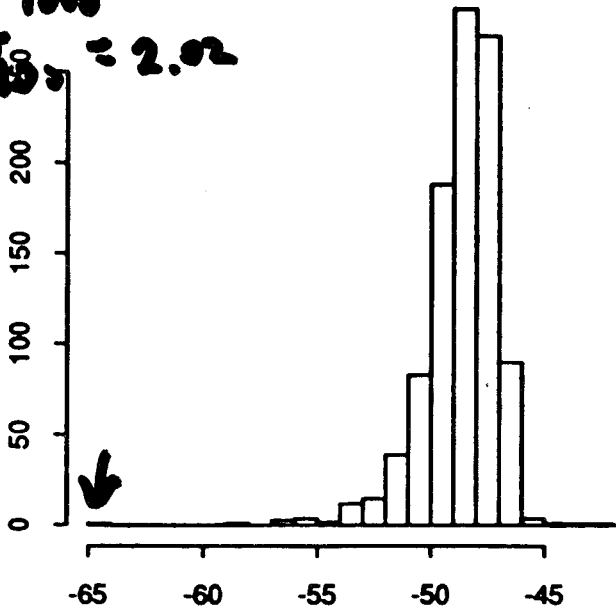
$$\hat{d} = 6074$$

# Bootstrap Analysis of $(\hat{a}, \hat{b}, \hat{c}, \hat{d})$

- Bootstrap sample  $\tilde{y}^* = (x_1^*, x_2^* \dots x_{51}^*)$   
a random sample, with replacement,  
of 51 pairs from  $\tilde{y} = (x_1, x_2, \dots x_{51})$ .
- Best replication of  $(\hat{a}, \hat{b}, \hat{c}, \hat{d})$  is  
 $(\hat{a}^*, \hat{b}^*, \hat{c}^*, \hat{d}^*)$   
the non-linear least-squares fit to  
the bootstrap sample  $\min_{\hat{a}} \sum_{i=1}^{51} [\hat{z}_i^* - \hat{z}_{0i}^*]^2$
- Did  $B=100$  and  $B=1000$  Bootstraps

1000 simple bootreps of ahat

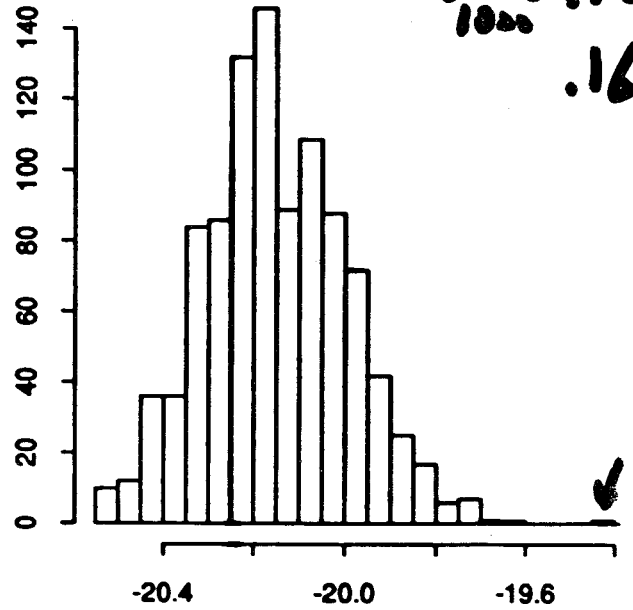
$\hat{\sigma}_{1000} = 1.71$   
 $\hat{\sigma}_{50} = 2.02$



aa

1000 simple bootreps of bhat

$\hat{\sigma}_{1000} = .157$   
.165

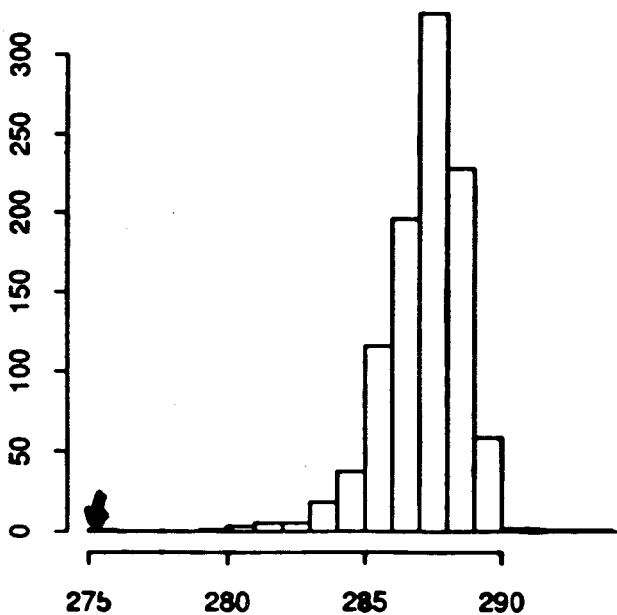


bb

NOT  
BELL-SHAPED!

$\hat{\sigma}_{1000} = 1.52$   
1.77

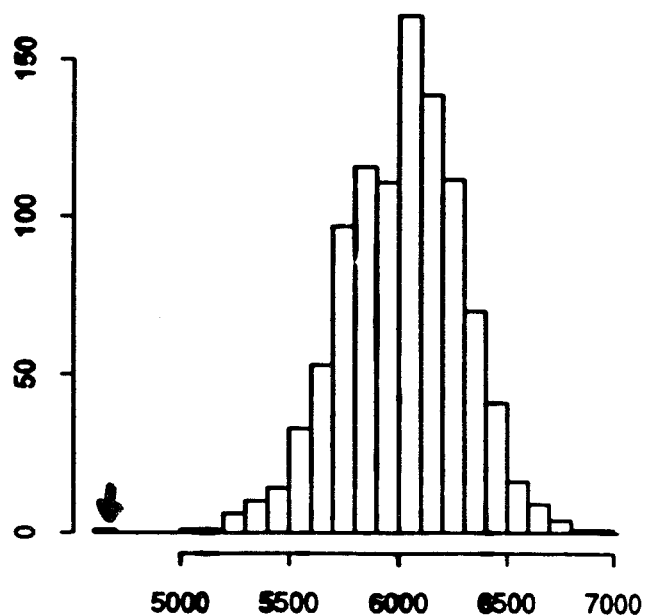
1000 simple bootreps of chat



cc

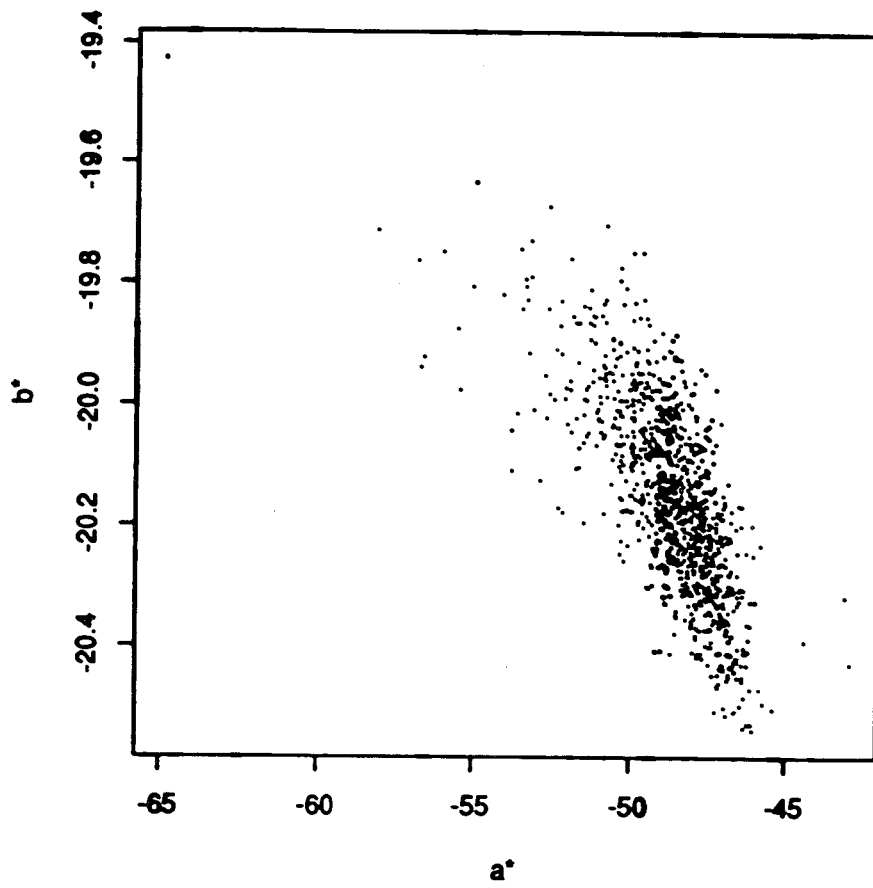
$\hat{\sigma}_{1000} = 277$   
308

1000 simple bootreps of dhat

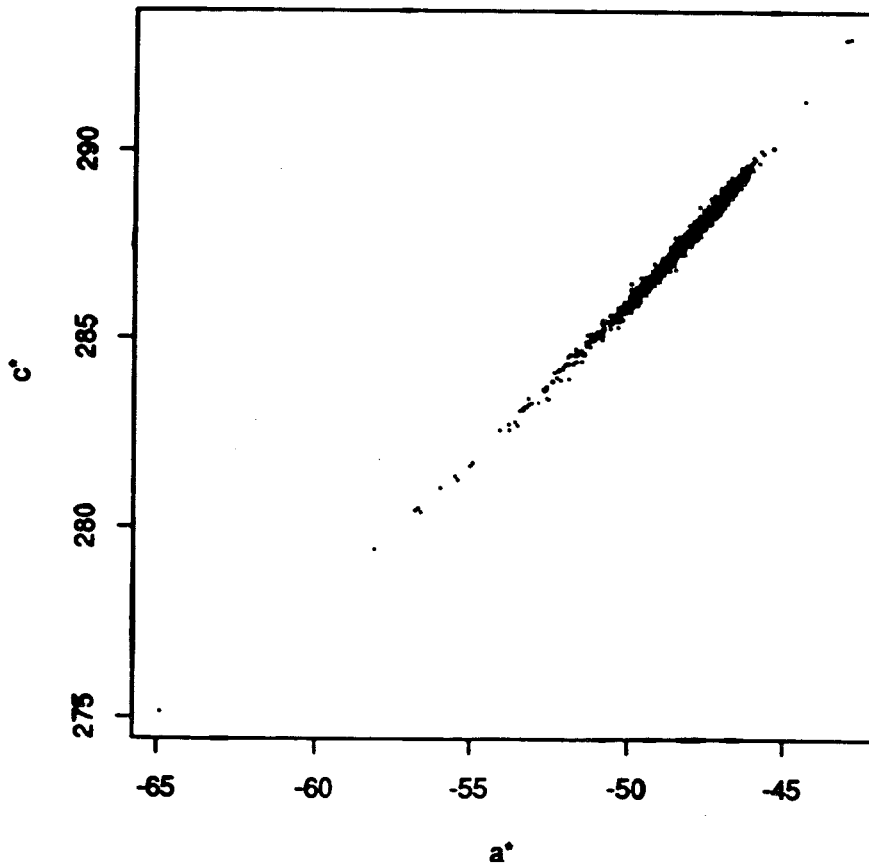


dd

1000 bootstrap pairs ( $a^*$ ,  $b^*$ )



1000 bootstrap pairs ( $a^*$ ,  $c^*$ )



# Non parametric Regression: "Lowess"

•  $n=164$  men,  $x_i = (\text{compliance}_i; \text{Chol. Decrease}_i)$

• Bootstrap:

$$\tilde{y} = (x_1^*, x_2^*, \dots, x_{164}^*)$$

• Lowess\*

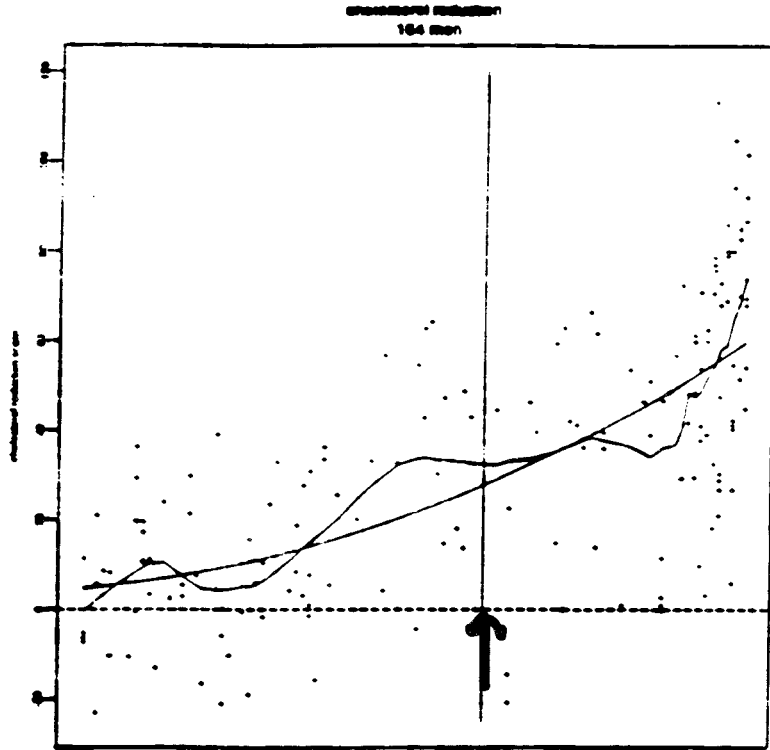
•  $B=50$

•  $\hat{\sigma}(\text{Lowess } 60\%)$

$= 5.71$

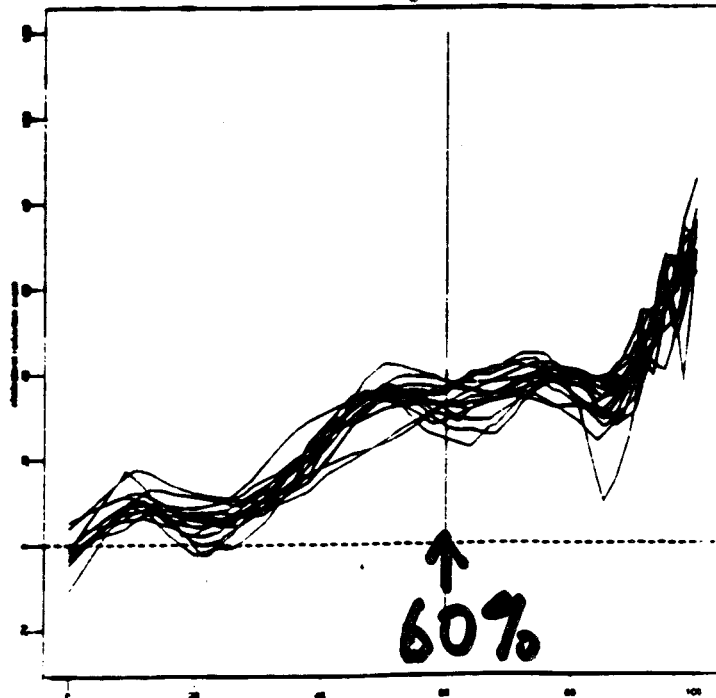
•  $\hat{\sigma}(\text{Bvad})$

$= 3.08$



COMPLIANCE →

cholesterol decrease →



60%

compliance

"The history of science exhibits a steady tendency to eliminate intellectual effort in the solution of individual problems, by developing comprehensive formulas which can resolve by rote a whole class of them."

... Ernest Nagel, 1955