

Statistical Theory

In the Age of the Computer:

BOOTSTRAP METHODS

Bradley Efron
(Stanford)

$n = 15$ 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
 $\bar{x} = .804$

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

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

$\hat{\sigma} = ? ?$

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

We have

DATA \underline{y}

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

PARAMETER OF Interest Θ

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

We Want to estimate

Θ from \underline{y}

Two Basic Questions

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

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

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

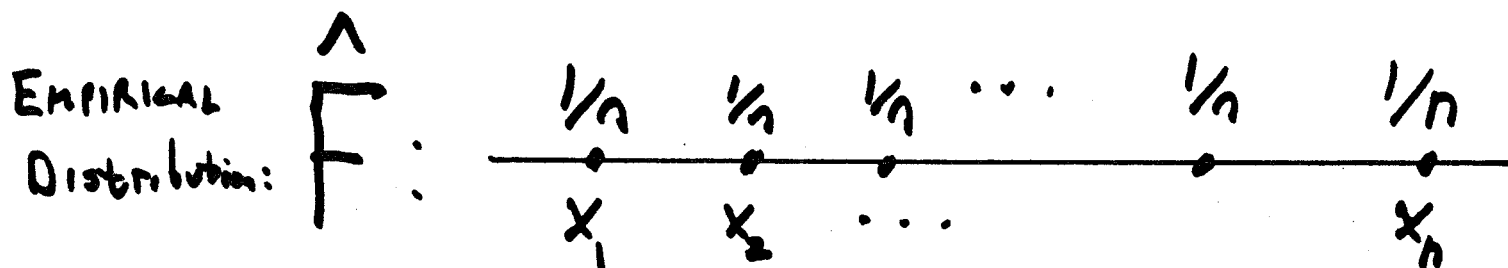
Statistic of Interest: $\hat{\theta}(y) = \bar{x} = \sum_{i=1}^n x_i / n$

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

Where

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

Estimating $\sigma(F)$



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

Gives estimated standard error

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

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

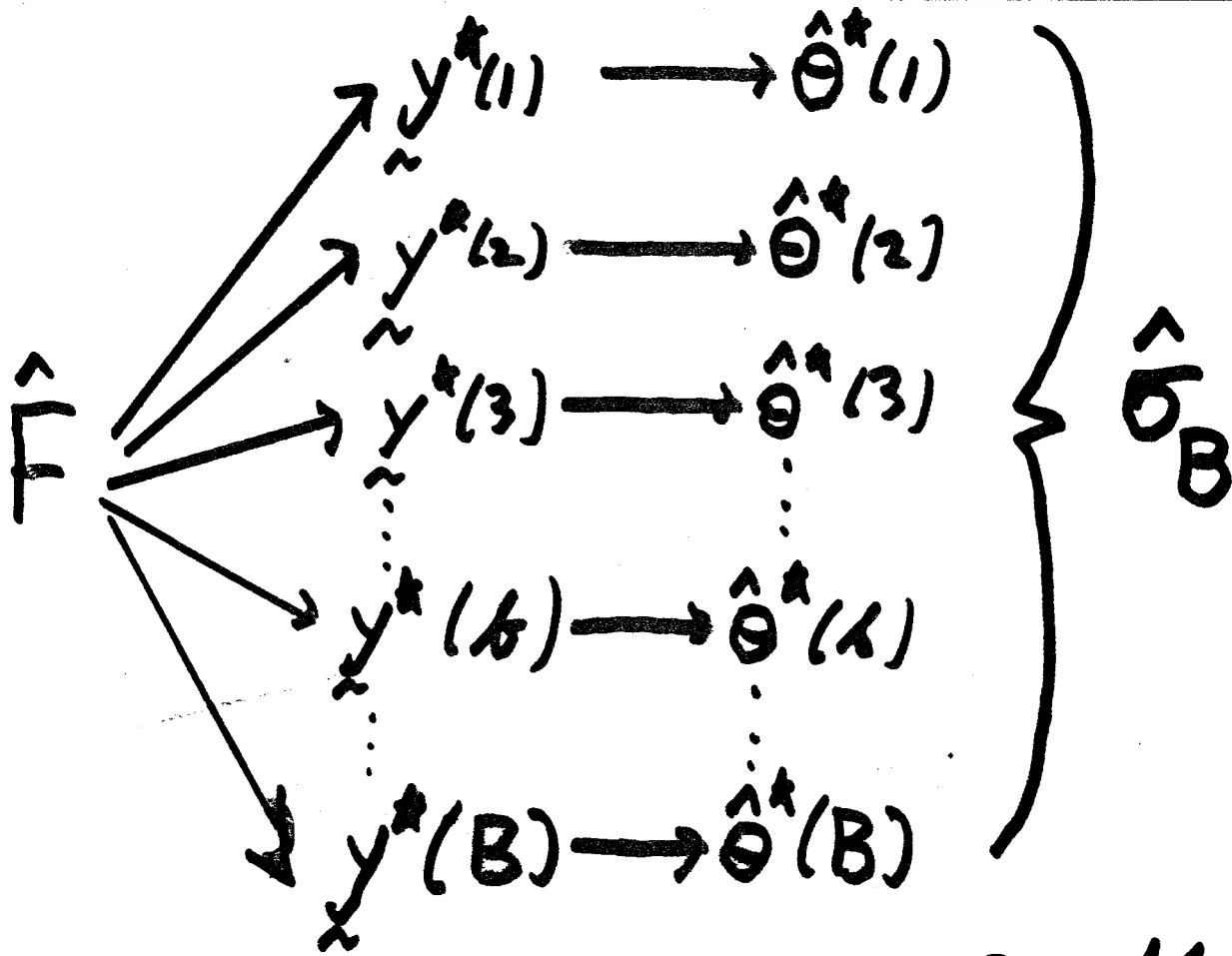
"Draw a random sample
of size n from \hat{F} "

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



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

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

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
------------------	-----	-----	-----	-----	-----

Results for the 15 lifetimes

B = 100 Bootstraps

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

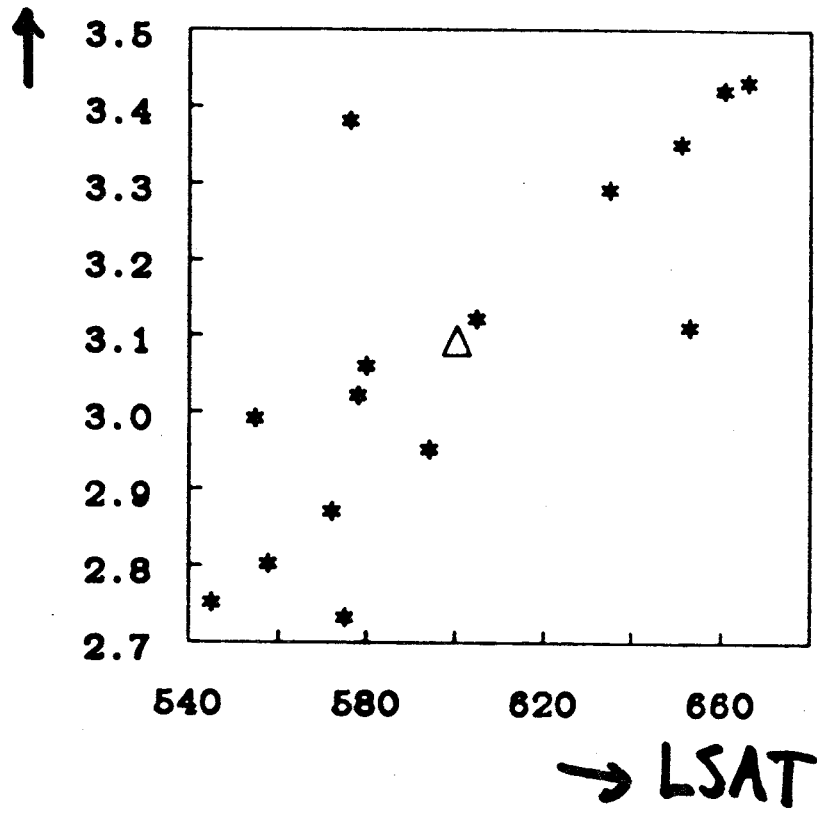
Sample Median $\hat{\sigma}_{100} = .229$

"boot(x, 100, "median")"

$$n = 15 \quad y = (x_1, x_2, \dots, x_{15})$$

GPA

15 LAW SCHOOLS



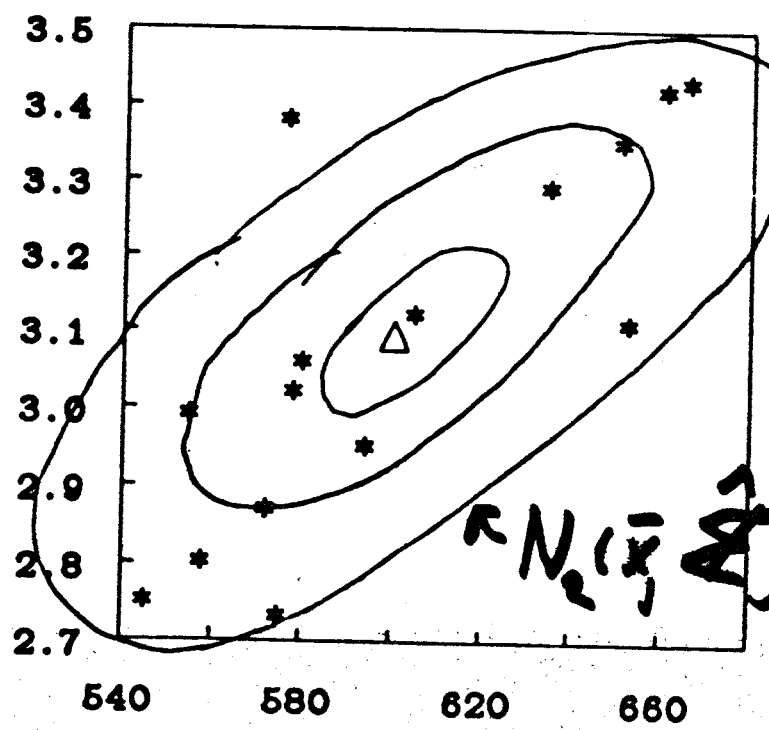
$$x_i = (\text{LSAT}_i, \text{GPA}_i)$$

$\hat{\theta}(y) =$ Correlation Coefficient

$$= .776$$

$$\hat{\sigma}_{1000} = .127$$

15 LAW SCHOOLS

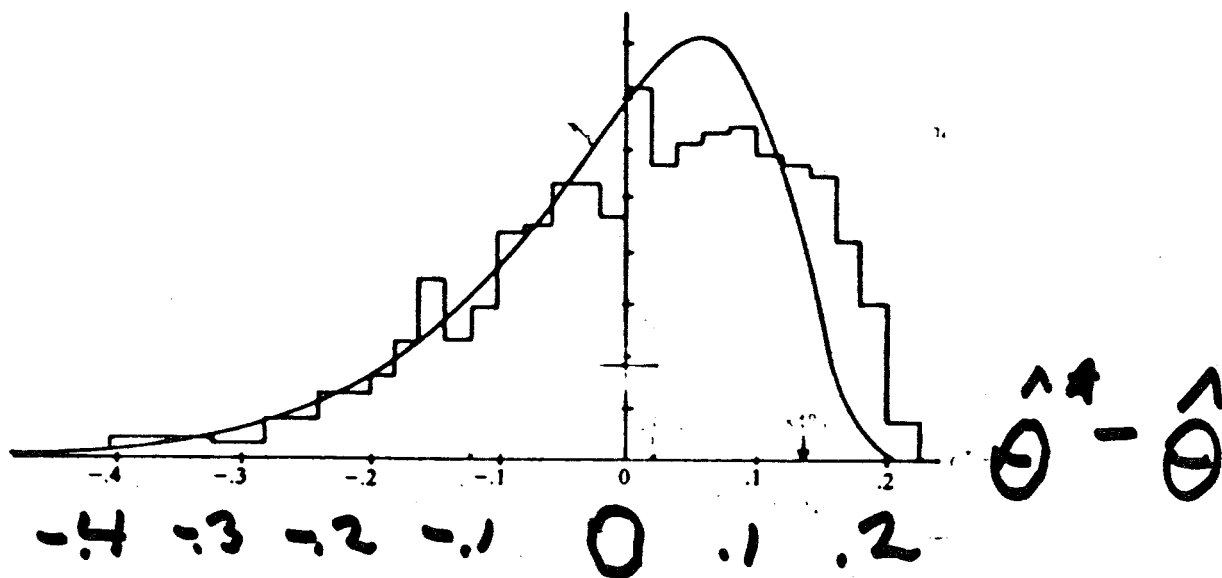


$$\hat{\sigma}_{\text{NORM}} = \frac{1 - \hat{\theta}^2}{\sqrt{12}}$$

$$= .115$$

$$N_2(\bar{x}, \hat{\Sigma}) = \hat{F}_{\text{NORM}}$$

Histogram of $B=1000 \hat{\theta}^*$



• $\hat{\sigma}_{1000} = .127$ compared to $\hat{\sigma}_{NORM} = 3.115$

• Smooth Curve is $B = \infty$ bootstraps starting from $\hat{F}_{NORM} \sim N_2(\bar{x}, \hat{\Sigma})$

$$\hat{F}_{NORM} \rightarrow \underset{\sim}{y}^* \rightarrow \hat{\theta}(\underset{\sim}{y}^*)$$

• Both histograms are very skewed to the left

A More Complicated Statistic

$n=88$ students each took five tests:

	Test 1	Test 2	Test 3	Test 4	Test 5	
Student 1	77	82	67	67	81	x_1
Student 2	44	56	55	61	36	x_2
Student 3	17	51	52	35	31	x_3
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Student 88	0	40	21	9	14	x_n
	Mech. (Closed)	Vectors (<)	Alg. (open)	Analysis (0)	Stats (0)	

$$X_{88 \times 5} \rightarrow X_c \rightarrow G_{5 \times 5} = X_c' X_c \rightarrow \text{eigenvectors}$$

1st eigenvector: (.51, .37, .35, .45, .53)

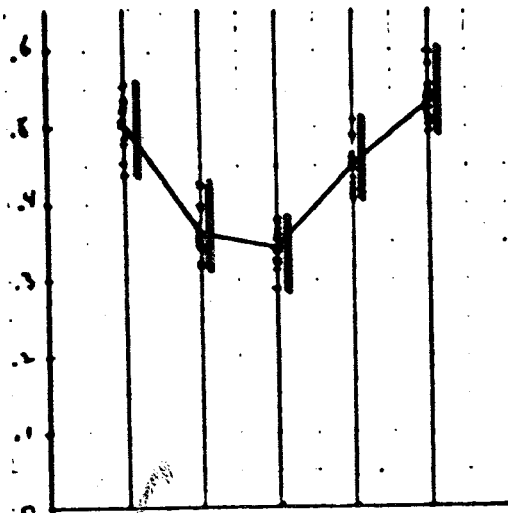
2nd eigenvector: (.75, .21, -.08, -.30, -.55)

} How
} Variable?

B=10 Bootstraps of the Test Data

First Eigenvector Components

1st 2nd 3rd 4th 5th



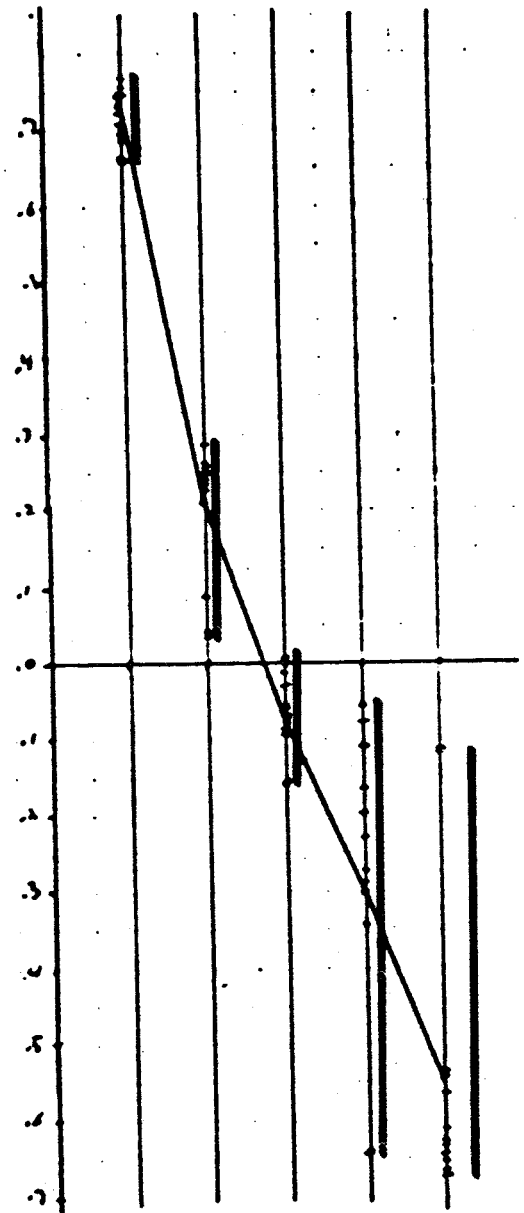
↪ 1st Eigenvector

Values of bootstrap eigenvalue components indicated by dashes.

2nd Eigenvector →

Second Eigenvector Components

1st 2nd 3rd 4th 5th



- Blue lines indicate range of boot values
- 2nd More variable than 1st

"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

References:

Efron and Tibshirani "Bootstrap Methods."

Statistical Science (86) Vol 3, No 3, p 54-77

Efron "Computers and the theory
of statistics..." SIAM REVIEW

(79) Vol 21, No 4 p 460-470