

Bootstrap Estimation

Suppose a simple random sample
(sampling with replacement)

$$X_1, X_2, \dots, X_n$$

is available from some population
with distribution function (cdf)

$$F(x).$$

Objective: Make inferences about
some feature of the population

- median
- variance
- correlation

1081

A statistic t_n is computed from the
observed data:

Sample mean: $t_n = \frac{1}{n} \sum_{j=1}^n X_j$

Standard deviation:

$$t_n = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X})^2}$$

Correlation:

$$X_j = \begin{bmatrix} X_{1j} \\ X_{2j} \end{bmatrix} \quad j = 1, 2, \dots, n$$

$$t_n = \frac{\sum_{j=1}^n (X_{1j} - \bar{X}_{1.})(X_{2j} - \bar{X}_{2.})}{\sqrt{\sum_{j=1}^n (X_{1j} - \bar{X}_{1.})^2} \sqrt{\sum_{j=1}^n (X_{2j} - \bar{X}_{2.})^2}}$$

1082

t_n estimates some feature of the
population.

What can you say about the dis-
tribution of t_n , with respect to all
possible samples of size n from the
population?

- Expectation of t_n
- Standard deviation
- Distribution function

How can a confidence interval be
constructed?

1083

Simulation: (The population c.d.f.
is known)

- For a univariate normal distribu-
tion with mean μ and variance
 σ^2 , the cdf is

$$\begin{aligned} F(x) &= Pr\{X \leq x\} \\ &= \int_{-\infty}^x \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{1}{2}\left(\frac{w-\mu}{\sigma}\right)^2} dw \end{aligned}$$

- Simulate B samples of size n and
compute the value of t_n for each
sample:

$$t_{n,1}, t_{n,2}, \dots, t_{n,B}$$

1084

- Approximate $E_F(t_n)$ with the simulated mean

$$\bar{t} = \frac{1}{B} \sum_{k=1}^B t_{n,k}$$

- Approximate $Var_F(t_n)$ with

$$\frac{1}{B-1} \sum_{k=1}^B (t_{n,k} - \bar{t})^2$$

- Approximate the standard deviation of t_n with

$$\sqrt{\frac{1}{B-1} \sum_{k=1}^B (t_{n,k} - \bar{t})^2}$$

1085

- Approximate the c.d.f. for t_n with

$$\bar{F}_n(t) = \frac{\text{number of samples with } t_{n,k} < t}{B}$$

- Order the B values of t_n from smallest to largest

$$t_{n(1)} \leq t_{n(2)} \leq \dots \leq t_{n(B)}$$

and approximate percentiles of the distribution of t_n .

1086

What if $F(X)$, the population c.d.f. is unknown?

- You cannot use a random number generator to simulate samples of size n and values of t_n , from the actual population.
- Use a bootstrap (or resampling) method?

1087

Basic idea:

- (1) Approximate the population cdf

$$F(x)$$

with the empirical cdf

$$\hat{F}_n(x)$$

obtained from the observed sample

$$X_1, X_2, \dots, X_n$$

Assign probability $\frac{1}{n}$ to each observation in the sample. Then

$$\bar{F}_n(x) = \frac{b}{n}$$

where b is the number of observations in the sample with

$$X_i < x \quad \text{for } i = 1, 2, \dots, n$$

1088

- **Approximate** the act of simulating B samples of size n from a population with c.d.f. $F(x)$ by simulating B samples of size n from a population with c.d.f. $\widehat{F}_n(x)$

- The “approximating” population is the original sample.
- Sample n observations from the original sample using simple random sampling with replacement.

This will be called a bootstrap sample.

1089

- repeat this B times to obtain B bootstrap samples of size n .

- Evaluate the summary statistic for each bootstrap sample

Sample 1: $t_{n,1}^*$

Sample 2: $t_{n,2}^*$

:

Sample B: $t_{n,B}^*$

1090

- Evaluate bootstrap estimates of features of the sampling distribution for t_n , when the c.d.f. is $\widehat{F}_n(x)$.

$$E_{\widehat{F}_n}(t_n) = \frac{1}{B} \sum_{b=1}^B t_{n,b}^*$$

$$Var_{\widehat{F}_n}(t_n)$$

$$= \frac{1}{B-1} \sum_{b=1}^B [t_{n,b}^* - E_{\widehat{F}_n}(t_n)]^2$$

1091

- This resampling procedure is called a *nonparametric bootstrap*

- If used properly it provides consistent large sample results:

As $n \rightarrow \infty$ and $B \rightarrow \infty$

$$E_{\widehat{F}_n}(t_n) = \frac{1}{B} \sum_{b=1}^B t_{n,b}^*$$

$$\rightarrow E_F(t_n)$$

$$Var_{\widehat{F}_n}(t_n) = \frac{1}{B-1} \sum_{b=1}^B \left[t_{n,b}^* - \frac{1}{B} \sum_{b=1}^B t_{n,b}^* \right]^2$$

$$\rightarrow Var_F(t_n)$$

1092

The bootstrap is a *large sample* method

- Large number of bootstrap samples. As $B \rightarrow \infty$

$$E_{\hat{F}_n}(t_n) = \frac{1}{B} \sum_{b=1}^B t_{n,b}^* \rightarrow E_{\hat{F}_n}(t_n)$$

$$\begin{aligned} \widehat{Var}_{\hat{F}_n}(t_n) &= \frac{1}{B-1} \sum_{b=1}^B (t_{n,b}^* - E_{\hat{F}_n}(t_n))^2 \\ &\rightarrow Var_{\hat{F}_n}(t_n) \end{aligned}$$

What is a good value for B ?

- standard deviation: $B \approx 200$
- confidence interval: $B \approx 1000$
- more demanding applications: $B \approx 5000$

1093

- Consistency: Original sample size must become large

As $n \rightarrow \infty$,

$$\hat{F}_n(x) \rightarrow F(x) \text{ for any } x$$

Then,

$$\begin{aligned} E_{\hat{F}_n}(t_n) &\rightarrow E_F(t_n) \\ Var_{\hat{F}_n}(t_n) &\rightarrow Var_F(t_n) \\ &\vdots \end{aligned}$$

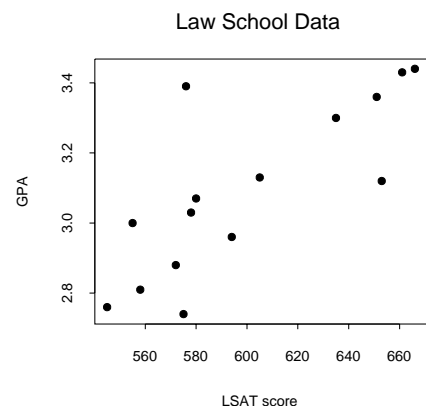
- For small values of n , $\hat{F}_n(x)$ could deviate substantially from $F(x)$.

1094

Example 12.1: Average values for GPA and LSAT scores for students admitted to $n=15$ Law Schools in 1973.

School	LSAT	GPA
1	576	3.39
2	635	3.30
3	558	2.81
4	579	3.03
5	666	3.44
6	580	3.07
7	555	3.00
8	661	3.43
9	651	3.36
10	605	3.13
11	653	3.12
12	575	2.74
13	545	2.76
14	572	2.88
15	594	2.96

1095



1096

- These schools were randomly selected from a larger population of law schools.
- We want to make inferences about the correlation between GPA and LSAT scores for the population of law schools
- The sample correlation ($n=15$) is

$$r = 0.7764 \equiv t_n$$

1097

Bootstrap samples

Take samples of $n=15$ schools, using simple random sampling with replacement

Sample 1:

School	LSAT	GPA
7	555	3.00
10	605	3.13
14	572	2.88
8	661	3.43
5	666	3.44
4	578	3.03
1	576	3.39
2	635	3.30
13	545	2.76
3	558	2.81
15	594	2.96
4	573	3.03
9	651	3.36
3	558	2.81
3	558	2.81

$$t_{n,1}^* = 0.8586 = r_{15,1}$$

1098

Sample 2:

School	LSAT	GPA
12	575	2.74
11	653	3.12
12	575	2.74
8	661	3.43
9	651	3.36
5	666	3.44
9	651	3.36
1	576	3.39
5	666	3.44
5	666	3.44
1	576	3.39
2	635	3.30
10	605	3.13
14	572	2.88

$$t_{n,2}^* = 0.6673 = r_{15,2}$$

1099

Repeat this $B = 5000$ times to obtain

$$t_{n,1}^*, t_{n,2}^*, \dots, t_{n,5000}^*$$

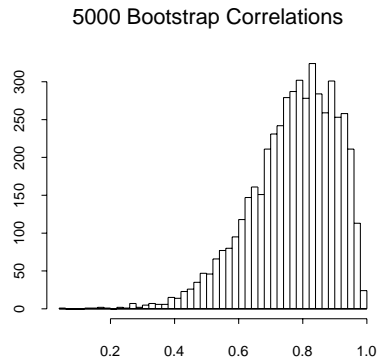
1100

Estimated correlation from the original sample of $n = 15$ law schools

$$r = 0.7764$$

Bootstrap standard error (from $B = 5000$ bootstrap samples) is

$$\begin{aligned} S_r &= \sqrt{\frac{1}{B-1} \sum_{b=1}^B \left[t_{n,b}^* - \frac{1}{B} \sum_{j=1}^B t_{n,j}^* \right]^2} \\ &= \sqrt{\frac{1}{4999} \sum_{b=1}^{5000} \left[r_{15,b} - \frac{1}{5000} \sum_{j=1}^{5000} r_{15,j} \right]^2} \\ &= 0.1341 \end{aligned}$$



1101

1102

Number of Bootstrap samples	Bootstrap estimate of Standard error
-----------------------------	--------------------------------------

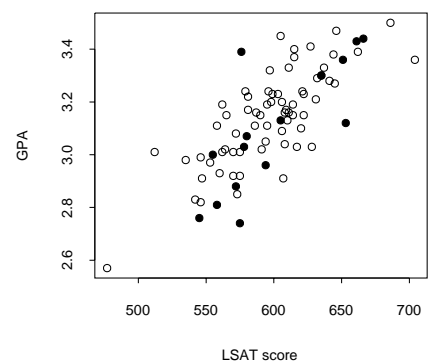
$B = 25$	0.1108
$B = 50$	0.0985
$B = 100$	0.1334
$B = 250$	0.1306
$B = 500$	0.1328
$B = 1000$	0.1299
$B = 2500$	0.1366
$B = 5000$	0.1341

“Very seldom are more than $B = 200$ replications needed for estimating a standard error.”

– Efron & Tibshirani (1993)
(page 52)

1103

Law School Data



1104

	School	LSAT	GPA	Type
1	1	622	3.23	2
2	2	542	2.83	2
3	3	579	3.24	2
4	4	653	3.12	1
5	5	606	3.09	2
6	6	576	3.39	1
7	7	620	3.10	2
8	8	615	3.40	2
9	9	553	2.97	2
10	10	607	2.91	2
11	11	558	3.11	2
12	12	596	3.24	2
13	13	635	3.30	1
14	14	581	3.22	2
15	15	661	3.43	1
16	16	547	2.91	2
17	17	599	3.23	2
18	18	646	3.47	2
19	19	622	3.15	2
20	20	611	3.33	2
21	21	546	2.99	2

1105

	School	LSAT	GPA	Type
22	22	614	3.19	2
23	23	628	3.03	2
24	24	575	3.01	2
25	25	662	3.39	2
26	26	627	3.41	2
27	27	608	3.04	2
28	28	632	3.29	2
29	29	587	3.16	2
30	30	581	3.17	2
31	31	605	3.13	1
32	32	704	3.36	2
33	33	477	2.57	2
34	34	591	3.02	2
35	35	578	3.03	1
36	36	572	2.88	1
37	37	615	3.37	2
38	38	606	3.20	2
39	39	603	3.23	2
40	40	535	2.98	2
41	41	595	3.11	2
42	42	575	2.92	2

1106

	School	LSAT	GPA	Type
43	43	573	2.85	2
44	44	644	3.38	2
45	45	545	2.76	1
46	46	645	3.27	2
47	47	651	3.36	1
48	48	562	3.19	2
49	49	609	3.17	2
50	50	555	3.00	1
51	51	586	3.11	2
52	52	580	3.07	1
53	53	594	2.96	1
54	54	594	3.05	2
55	55	560	2.93	2
56	56	641	3.28	2
57	57	512	3.01	2
58	58	631	3.21	2
59	59	597	3.32	2
60	60	621	3.24	2
61	61	617	3.03	2
62	62	637	3.33	2
63	63	572	3.08	2

1107

	School	LSAT	GPA	Type
64	64	610	3.13	2
65	65	562	3.01	2
66	66	635	3.30	2
67	67	614	3.15	2
68	68	546	2.82	2
69	69	598	3.20	2
70	70	666	3.44	1
71	71	570	3.01	2
72	72	570	2.92	2
73	73	605	3.45	2
74	74	565	3.15	2
75	75	686	3.50	2
76	76	608	3.16	2
77	77	595	3.19	2
78	78	590	3.15	2
79	79	558	2.81	1
80	80	611	3.16	2
81	81	564	3.02	2
82	82	575	2.74	1

1108

The correlation coefficient for the population of 82 Law Schools in 1973 is

$$\rho = 0.7600$$

The exact distribution of estimated correlation coefficients for random samples of size $n=15$ from this population involves

More than 3.6×10^{110} possible samples of size $n=15$

1109

Results from 100,000 samples of size $n=15$.

Percentiles	
.95	0.9128
.75	0.8404
.50	0.7699
.25	0.6777
.05	0.5010

$$\text{min} = -0.3647$$

$$\text{max} = 0.9856$$

$$\text{mean} = 0.7464$$

$$\text{std. error} = 0.1302$$

1110

Bias: From a sample of size n , t_n is used to estimate a population parameter θ .

$$Bias_F(t_n) = E_F(t_n) - \theta$$

\uparrow \uparrow
 average across true
 all possible parameter
 samples of size n value
 from the population

Law School example:

$$\begin{aligned}
 Bias_F(r_{15}) &= E_F(r_{15}) - 0.7600 \\
 &= .7464 - .7600 \\
 &= -0.0136
 \end{aligned}$$

1111

Bootstrap estimate of bias:

“true value” if you take the original sample as the population

↓

$$Bias_{\bar{F}}(t_n) = E_{\bar{F}}(t_n) - t_n$$

↑

approximate this with the average of results from B bootstrap samples $\frac{1}{B} \sum_{b=1}^B t_{n,b}^*$

1112

Law School example:

($B = 5000$ bootstrap samples)

$$\begin{aligned}\widehat{Bias}_F(r_{15}) &= \frac{1}{5000} \sum_{b=1}^{5000} r_{15,b} - r_{15} \\ &= .7692 - .7764 \\ &= -.0072\end{aligned}$$

1113

Improved bootstrap bias estimation: Efron & Tibshirani (1993), Section 10.4

Bias corrected estimates:

$$\begin{aligned}\bar{t}_n &= t_n - \widehat{Bias}_b(t_n) \\ &= 2t_n - \left(\frac{1}{B} \sum_{b=1}^B t_{n,b}^* \right)\end{aligned}$$

Law School example:

$$\begin{aligned}\tilde{r}_{15} &= r_{15} - \widehat{Bias}_b(r_{15}) \\ &= .7764 - (-.0072) = 0.7836\end{aligned}$$

↑
here we moved
farther away from
 $\rho = .7600$.

1114

Bias correction can be dangerous in practice:

- \bar{t}_n may have a substantially larger variance than t_n
- $MSE_F(\bar{t}_n)$
 $= [Bias_F(\bar{t}_n)]^2 + Var_F(\bar{t}_n)$

is often larger than

$$\begin{aligned}MSE_F(t_n) \\ = [Bias_F(t_n)]^2 + Var_F(t_n)\end{aligned}$$

1115

Empirical Percentile Bootstrap Confidence Intervals

Construct an approximate

$$(1 - \alpha) \times 100\%$$

confidence interval for θ .

- t_n is an estimator for θ
- Compute B bootstrap samples to obtain $t_{n,1}^*, \dots, t_{n,B}^*$
- Order the bootstrapped values from smallest to largest

$$t_{n,(1)} \leq t_{n,(2)} \leq \dots \leq t_{n,(B)}$$

1116

- Compute upper and lower $\frac{\alpha}{2} \times 100$ th percentiles

Compute

$$k_L = \left\lceil (B + 1) \frac{\alpha}{2} \right\rceil$$

$$= \text{largest integer} \leq (B + 1) \frac{\alpha}{2}$$

$$k_U = B + 1 - k_L$$

Then, an approximate

$$(1 - \alpha) \times 100\%$$

confidence interval for θ is

$$[t_{n,(k_L)}, t_{n,(k_U)}]$$

1117

Law School example:

($B = 5000$ bootstrap samples)

Construct an approximate 90% confidence interval for

$\rho =$ population correlation

$$\alpha = .10$$

$$k_L = [(5001)(.05)] = [250.05] = 250$$

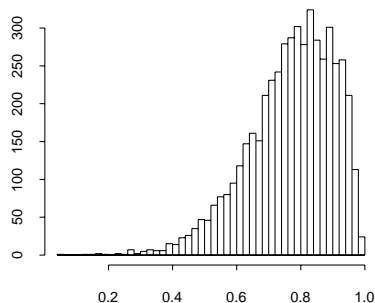
$$k_U = 5001 - 250 = 4751$$

An approximate large sample 90% confidence interval is

$$[r_{15,(250)}, r_{15,(4751)}] = [0.5207, 0.9487]$$

1118

5000 Bootstrap Correlations



1119

Law School example:

- Fisher Z -transformation

$$Z_n = \frac{1}{2} \log \left(\frac{1 + r_n}{1 - r_n} \right)$$

$$\sim N \left(\frac{1}{2} \log \left(\frac{1 + \rho}{1 - \rho} \right), \frac{1}{n - 3} \right)$$

- An approximate $(1 - \alpha) \times 100\%$ confidence interval for $\frac{1}{2} \log \left(\frac{1 + \rho}{1 - \rho} \right)$ is

$$\text{lower limit: } Z_n - Z_{\alpha/2} \frac{1}{\sqrt{n-3}} = Z_L$$

$$\text{upper limit: } Z_n + Z_{\alpha/2} \frac{1}{\sqrt{n-3}} = Z_U$$

- Transform back to the original scale

$$\left[\frac{e^{2Z_L} - 1}{e^{2Z_L} + 1}, \frac{e^{2Z_U} - 1}{e^{2Z_U} + 1} \right]$$

1120

For $n = 15$ and $r_{15} = .7764$ we have

$$Z_{15} = \frac{1}{2} \log \left(\frac{1 + .7764}{1 - .7764} \right) = 1.03624$$

and

$$Z_L = Z_{15} - Z_{.05} \frac{1}{\sqrt{15-3}} = .561372$$

$$Z_U = Z_{15} + Z_{.05} \frac{1}{\sqrt{15-3}} = 1.511113$$

and an approximate 90% confidence interval for the correlation is

$$(0.509, 0.907)$$

The bootstrap percentile interval would approximate this interval if the original sample was taken from a bivariate normal approximation.

1121

Percentile Bootstrap Confidence Intervals

Suppose there is a transformation

$$\phi = m(\theta)$$

such that

$$\bar{\phi} = m(t_n) \sim N(\phi, \omega^2)$$

for some standard deviation ω . Then, an approximate confidence interval for θ is

$$[m^{-1}(\bar{\phi} - Z_{(\alpha/2)}\omega), m^{-1}(\bar{\phi} + Z_{\alpha/2}\omega)]$$

1122

- The bootstrap percentile interval is a consistent approximation. (You do not have to identify the $m()$ transformation.)
- The bootstrap approximation becomes more accurate for larger sample
- For smaller samples, the coverage probability of the bootstrap percentile interval tends to be smaller than the nominal level $(1 - \alpha) \times 100$.

1123

- A percentile interval is entirely inside the parameter space.

A percentile confidence interval for a correlation lies inside the interval $[-1, 1]$.

1124

Bias-corrected and accelerated (BC_a) bootstrap percentile confidence intervals

- simulate B bootstrap samples and order the resulting estimates

$$t_{n,(1)}^* \leq t_{n,(2)}^* \leq \dots \leq t_{n,(B)}^*$$

1125

- The BC_a interval of intended coverage $1 - \alpha$ is given by

$$[t_{n,([\alpha_1(B+1)])}^*, t_{n,([\alpha_2(B+1)])}^*]$$

where

$$\alpha_1 = \Phi \left(\bar{Z}_0 + \frac{\bar{Z}_0 - Z_{\alpha/2}}{1 - \bar{a}(\bar{Z}_0 - Z_{\alpha/2})} \right)$$

$$\alpha_2 = \Phi \left(\bar{Z}_0 + \frac{\bar{Z}_0 + Z_{\alpha/2}}{1 - \bar{a}(\bar{Z}_0 + Z_{\alpha/2})} \right)$$

1126

and

$\Phi()$ is the c.d.f. for the standard normal distribution

$Z_{\alpha/2}$ is an “upper” percentile of the standard normal distribution, e.g., $Z_{.05} = 1.645$ and $\Phi(Z_{\alpha/2}) = 1 - \frac{\alpha}{2}$

$$\bar{Z}_0 = \Phi^{-1} \left(\frac{\text{proportion of } t_{n,b}^* \text{ values smaller than } t_n}{1} \right)$$

is roughly a measure of median bias of t_n in “normal units.”

When exactly half of the bootstrap samples have $t_{n,b}^*$ values less than t_n , then $\bar{Z}_0 = 0$.

1127

$$\bar{a} = \frac{\sum_{j=1}^n (t_{n,(.)} - t_{n,-j})^3}{6 \left[\sum_{j=1}^n (t_{n,(.)} - t_{n,-j})^2 \right]^{3/2}}$$

is the estimated acceleration, where

$t_{n,-j}$ is the value of t_n when the j -th case is removed from the sample

and

$$t_{n,(.)} = \frac{1}{n} \sum_{j=1}^n (t_{n,-j})$$

1128

- BC_a intervals are second order accurate

$Pr\{\theta < \text{lower end of } BC_a \text{ interval}\}$

$$= \frac{\alpha}{2} + \frac{C_{lower}}{n}$$

$Pr\{\theta > \text{upper end of } BC_a \text{ interval}\}$

$$= \frac{\alpha}{2} + \frac{C_{upper}}{n}$$

- Bootstrap percentile intervals are first order accurate

$$Pr\{\theta < \text{lower end}\} = \frac{\alpha}{2} + \frac{C_{lower}^*}{\sqrt{n}}$$

$$Pr\{\theta < \text{upper end}\} = \frac{\alpha}{2} + \frac{C_{upper}^*}{\sqrt{n}}$$

1129

- ABC intervals are approximations to BC_a intervals

- second order accurate
- only use about 3% of the computation time

(See Efron & Tibshirani (1993) Chapter 14)

1130

```
# This is S-plus code for creating
# bootstrapped confidence intervals
# for a correlation coefficient. It is
# stored in the file
#
#         lawschl.ssc
#
# Any line preceded with a pound sign
# is a comment that is ignored by the
# program. The law school data are
# read from the file
#
#         lawschl.dat
#
# Enter the law school data into a data frame

laws <- read.table("lawschl.dat",
                  col.names=c("School", "LSAT", "GPA"))
laws
```

1131

	School	LSAT	GPA
1	1	576	3.39
2	2	635	3.30
3	3	558	2.81
4	4	578	3.03
5	5	666	3.44
6	6	580	3.07
7	7	555	3.00
8	8	661	3.43
9	9	651	3.36
10	10	605	3.13
11	11	653	3.12
12	12	575	2.74
13	13	545	2.76
14	14	572	2.88
15	15	594	2.96

1132

```

# Plot the data

par(fin=c(7.0,7.0),pch=16,mkh=.15,mex=1.5)
plot(laws$LSAT,laws$GPA, type="p",
      xlab="GPA",ylab="LSAT score",
      main="Law School Data")

# Compute the sample correlation matrix

rr<-cor(laws$LSAT,laws$GPA)
cat("Estimated correlation: ",
    round(rr,5), fill=T)

Estimated correlation:  0.77637

```

1133

```

# First test for zero correlation

n <- length(laws$LSAT);
tt<- sqrt(n-2)*rr/sqrt(1 - rr*rr)
pval <- 1 - pt(tt,n-2)
pval <- round(pval,digits=5)

cat("t-test for zero correlation: ",
    round(tt,4), fill=T)

t-test for zero correlation:  4.4413

cat("p-values for the t-test for
zero correlation: ", pval, fill=T)

p-values for the t-test
for zero correlation:  0.00033

```

1134

```

# Use Fisher's z-transformation to construct
# approximate confidence intervals.
# First set the level of confidence at 1-alpha.

alpha <- .10
z <- 0.5*log((1+rr)/(1-rr))
z1 <- z - qnorm(1-alpha/2)/sqrt(n-3)
zu <- z + qnorm(1-alpha/2)/sqrt(n-3)
r1 <- round((exp(2*z1)-1)/(exp(2*z1)+1),
            digits=4)
ru <- round((exp(2*zu)-1)/(exp(2*zu)+1),
            digits=4)

per <- (1-alpha)*100;
cat( per,"% confidence interval: (",
    r1,", ",ru,")",fill=T)

90 % confidence interval:  ( 0.509 ,  0.9071 )

```

1135

```

# Compute bootstrap confidence intervals.
# Use B=5000 bootstrap samples.

nboot <- 5000
rboot <- bootstrap(data=laws,
                   statistic=cor(GPA,LSAT),B=nboot)

Forming replications  1  to 100
Forming replications 101 to 200
Forming replications 201 to 300
.
.
.
Forming replications 4801 to 4900
Forming replications 4901 to 5000

```

1136

```

# limits.emp(): Calculates empirical percentiles
#               for the bootstrapped parameter
#               estimates in a resamp object.
#               The quantile function is used to
#               calculate the empirical percentiles.
# usage:
# limits.emp(x, probs=c(0.025, 0.05, 0.95, 0.975))

```

```

limits.emp(rboot, probs=c(0.05,0.95))

          5%      95%
Param 0.523511 0.9473424

```

```

# Do another set of 5000 bootstrapped values

```

```

rboot <- bootstrap(data=laws,
                   statistic=cor(GPA,LSAT),B=nboot)
limits.emp(rboot, probs=c(0.05,0.95))

```

```

          5%      95%
Param 0.520661 0.9486835

```

1137

```

# limits.bca(): Calculates BCa (bootstrap
#               bias-correct, adjusted)
#               confidence limits.
# usage:

```

```

# limits.bca(boot.obj,
#            probs=c(0.025, 0.05, 0.95, 0.975),
#            details=F, z0=NULL,
#            acceleration=NULL,
#            group.size=NULL,
#            frame.eval.jack=sys.parent(1))

```

```

limits.bca(rboot,probs=c(0.05,0.95),detail=T)

```

```

$limits:
          5%      95%
Param 0.442216 0.9306627

```

```

$emp.probs:
          5%      95%
Param 0.01954177 0.9065055

```

1138

```

$z0:
      Param
-0.07979538

```

```

$acceleration:
      Param
-0.07567156

```

```

$group.size:
[1] 1

```

1139

```

# Both sets of confidence intervals could
# have been obtained from the summary( )
# function

```

```

summary(rboot)

```

```

Call:
bootstrap(data = laws, statistic =
cor(GPA, LSAT), B = nboot)

```

```

Number of Replications: 5000

```

```

Summary Statistics:
      Observed   Bias   Mean   SE
Param  0.7764 -0.007183 0.7692 0.1341

```

```

Empirical Percentiles:
      2.5%   5%   95%  97.5%
Param 0.4638 0.5207 0.9487 0.9616

```

```

BCa Percentiles:
      2.5%   5%   95%  97.5%
Param 0.3462 0.4422 0.9307 0.9449

```

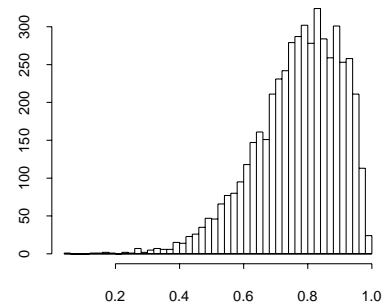
1140

```
# Make a histogram of the bootstrapped
# correlations.
```

```
hist(rboot$rep,nclass=50,
     xlab=" ",
     main="5000 Bootstrap Correlations",
     density=.0001)
```

1141

5000 Bootstrap Correlations



1142

The bootstrap can fail:

Example: X_1, X_2, \dots, X_n are sampled from a uniform $(0, \theta)$ distribution:

true density:

$$f(x) = \frac{1}{\theta}, \quad 0 < x < \theta$$

$$\text{true c.d.f.: } F(x) = \begin{cases} 0 & x \leq 0 \\ \frac{x}{\theta} & 0 < x \leq \theta \\ 1 & x > \theta \end{cases}$$

Bootstrap percentile confidence intervals for θ tend to be too short.

1143

Application of the bootstrap must adequately replicate the random process that produced the original sample

- Simple random samples
- “Nested” experiments
 - Sample plants from a field
 - Sample leaves from plants
- Curve fitting (existence of covariates)
 - Fixed levels
 - Random samples

1144

Parametric Bootstrap

- Suppose you “knew” that

$$(X_{1j}, X_{2j}) \quad j = 1, \dots, n$$

were obtained from a simple random sample from a bivariate normal distribution, i.e.,

$$X_j = \begin{bmatrix} X_{1j} \\ X_{2j} \end{bmatrix} \sim NID \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \Sigma \right)$$

- Estimate unknown parameters

$$\bar{\mu} = \begin{bmatrix} \bar{\mu}_1 \\ \bar{\mu}_2 \end{bmatrix} = \frac{1}{n} \sum_{j=1}^n X_j = \bar{X}$$

$$\bar{\Sigma} = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X})(X_j - \bar{X})^T$$

1145

- Obtain a bootstrap sample of size n by sampling from a $N(\hat{\mu}, \hat{\Sigma})$ distribution, i.e.

$$X_{1,b}, \dots, X_{n,b}$$

and compute $t_{n,b}^* = r_{n,b}$

- Repeat this to obtain

$$r_{n,1}, r_{n,2}, \dots, r_{n,B}$$

- Compute bootstrap

- standard errors
- bias estimators
- confidence intervals

1146

References

Davison, A.C. and Hinkley, D.V. (1997) **Bootstrap Methods and Their Applications**, Cambridge Series in Statistical and Probabilistic Mathematics 1, Cambridge University Press, New York.

Efron, B. (1982) **The Jackknife, The Bootstrap and other resampling plans**, CBMS, 38, SIAM-NSF, Philadelphia.

Efron, B. and Gong, G. (1983) **The American Statistician**, , 36-48.

Efron, B. (1987) **Better Bootstrap confidence intervals (with discussion)** *Journal of the American Statistical Association*, 82, 171-200.

Efron, B. and Tibsharani, R. (1993) **An Introduction to the Bootstrap**, Chapman and Hall, New York.

Shao, J. and Tu, D. (1995) **The Jackknife and Bootstrap**, New York, Springer.

1147

Example 12.2: Stormer viscometer data (*Venables & Ripley, Chapter 8*)

- measure viscosity of a fluid
- measure time taken for an inner cylinder in the mechanism to complete a specific number of revolutions in response to an actuating weight
- calibrate the viscometer using runs with
 - varying weights (W) (g)
 - fluids with known viscosity (V)
 - record the time (T) (sec)

1148

● theoretical model

$$T = \frac{\beta_1 V}{W - \beta_2} + \epsilon$$

```
# This code is used to explore the
# Stormer viscometer data. It is stored
# in the file
#
#           stormer.ssc
#
# Enter the data into a data frame.
# The data are stored in the file
#
#           stormer.dat

library(MASS)

stormer <- read.table("stormer.dat")
stormer
```

1149

	Viscosity	Wt	Time
1	14.7	20	35.6
2	27.5	20	54.3
3	42.0	20	75.6
4	75.7	20	121.2
5	89.7	20	150.8
6	146.6	20	229.0
7	158.3	20	270.0
8	14.7	50	17.6
9	27.5	50	24.3
10	42.0	50	31.4
11	75.7	50	47.2
12	89.7	50	58.3
13	146.6	50	85.6
14	158.3	50	101.1
15	161.1	50	92.2
16	298.3	50	187.2
17	75.7	100	24.6
18	89.7	100	30.0
19	146.6	100	41.7
20	158.3	100	50.3
21	161.1	100	45.1
22	298.3	100	89.0
23	298.3	100	86.5

1150

Starting values for β_1 and β_2 :

Fit an approximate linear model.
Note that

$$T_i = \frac{\beta_1 V_i}{W_i - \beta_2} + \epsilon_i$$

$$\Rightarrow (W_i - \beta_2)T_i = \beta_1 V_i + \epsilon_i(W_i - \beta_2)$$

$$\Rightarrow W_i T_i = \beta_1 V_i + \beta_2 T_i + \epsilon_i(W_i - \beta_2)$$

↑
↑

this is the new
this is the

response variable
new error

Use OLS estimation to obtain

$$\bar{\beta}_1^{(0)} = 28.876$$

$$\bar{\beta}_2^{(0)} = 2.8437$$

1151

```
# Use a linear approximation to obtain
# starting values for least squares
# estimation in the non-linear model

fm0 <- lm(Wt*Time ~ Viscosity + Time - 1,
          data=stormer)

b0 <- coef(fm0)
names(b0) <- c("b1","b2")

# Fit the non-linear model

storm.fm <- nls(
  formula = Time ~ b1*Viscosity/(Wt-b2),
  data = stormer, start = b0,
  trace = T)
```

```
885.365 : 28.8755 2.84373
825.11  : 29.3935 2.23328
825.051 : 29.4013 2.21823
```

1152

```
summary(storm.fm)$parameters
```

	Value	Std. Error	t value
b1	29.401328	0.9155353	32.113813
b2	2.218226	0.6655234	3.333054

1153

```
# Create a bivariate confidence region  
# for the the (b1,b2) parameters.  
# First set up a grid of (b1,b2) values
```

```
bc <- coef(storm.fm)  
se <- sqrt(diag(vcov(storm.fm)))  
dv <- deviance(storm.fm)
```

```
gsize<-51  
b1 <- bc[1] + seq(-3*se[1], 3*se[1],  
                 length = gsize)  
b2 <- bc[2] + seq(-3*se[2], 3*se[2],  
                 length = gsize)  
bv <- expand.grid(b1, b2)
```

```
# Create a function to evaluate sums of squares
```

```
ssq <- function(b)  
  sum((stormer$Time - b[1] * stormer$Viscosity/  
      (stormer$Wt-b[2]))^2)
```

1154

```
# Evaluate the sum of squared residuals and  
# approximate F-ratios for all of the  
# (b1,b2) values on the grid
```

```
dbeta <- apply(bv, 1, ssq)  
n<-length(stormer$Viscosity)  
df1<-length(bc)  
df2<-n-df1  
fstat <- matrix( ((dbeta - dv)/df1) / (dv/df2),  
                gsize, gsize)
```

1155

```
# Create the plot
```

```
par(fin=c(7.0,7.0), mex=1.5,lwd=3)  
plot(b1, b2, type="n",  
     main="95% Confidence Region")  
contour(b1, b2, fstat, levels=c(1,2,5,7,10,15,20),  
        labex=0.75, lty=2, add=T)  
contour(b1, b2, fstat, levels=qf(0.95,2,21),  
        labex=0, lty=1, add=T)  
text(31.6,0.3,"95% CR", adj=0, cex=0.75)  
points(bc[1], bc[2], pch=3, mkh=0.15)
```

```
# remove b1,b2, and bv
```

```
rm(b1,b2,bv,fstat)
```

1156

Construct a joint confidence region for (β_1, β_2) :

- Deviance (residual sum of squares):

$$d(\beta_1, \beta_2) = \sum_{j=1}^n \left[T_j - \frac{\beta_1 V_j}{W_j - \beta_2} \right]^2$$

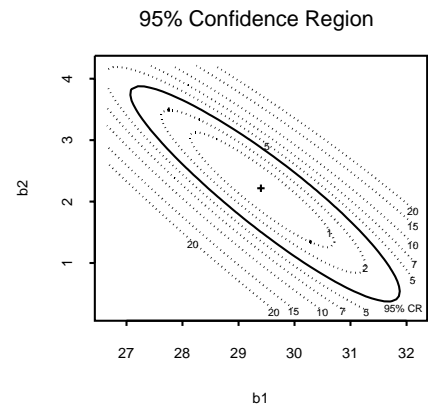
- Approximate F-statistic:

$$F(\beta_1, \beta_2) = \frac{\frac{d(\beta_1, \beta_2) - d(\hat{\beta}_1, \hat{\beta}_2)}{2}}{\frac{d(\hat{\beta}_1, \hat{\beta}_2)}{n-2}}$$

- An approximate $(1 - \alpha) \times 100\%$ confidence interval consists of all (β_1, β_2) such that

$$F(\beta_1, \beta_2) < F_{(2, n-2), \alpha}$$

1157



1158

Bootstrap Estimation

Bootstrap I:

Sample $n = 23$ cases

$$(T_i, W_i, V_i)$$

from the original data set
(using simple random
sampling with replacement)

1159

```
#####
# Bootstrap I:
#   Treat the regressors as random and
#   resample the cases (y,x_1,x_2)
#####

storm.boot1 <- bootstrap(stormer,
  coef(nls(Time~b1*Viscosity/(Wt-b2),
  data=stormer,start=bc)),B=1000)

summary(storm.boot1)

Call:
bootstrap(data = stormer,
  statistic = coef(nls(Time ~ (b1 *
  Viscosity)/(Wt - b2),
  data = stormer, start = bc)), B = 1000)

Number of Replications: 1000
```

1160

Summary Statistics:

	Observed	Bias	Mean	SE
b1	29.401	-0.08665	29.315	0.7194
b2	2.218	0.09412	2.312	0.8424

Empirical Percentiles:

	2.5%	5%	95%	97.5%
b1	27.8202	28.113	30.402	30.613
b2	0.8453	1.057	3.513	3.779

BCa Percentiles:

	2.5%	5%	95%	97.5%
b1	27.9368	28.242	30.525	30.747
b2	0.7462	0.899	3.409	3.604

Correlation of Replicates:

	b1	b2
b1	1.0000	-0.8628
b2	-0.8628	1.0000

1161

```
# Produce histograms of the bootstrapped
# values of the regression coefficients

# The following code is used to draw
# non-parametric estimated densities,
# normal densities, and a histogram
# on the same graph.

# truehist(): Plot a Histogram (prob=T
# by default) For the function
# hist(), probability=F by
# default.
# width.SJ(): Bandwidth Selection by Pilot
# Estimation of Derivatives.
# Uses the method of Sheather
# & Jones (1991) to select the
# bandwidth of a Gaussian kernel
# density estimator.
```

1162

```
# density() : Estimate Probability Density
# Function. Returns x and y
# coordinates of a non-parametric
# estimate of the probability
# density of the data. Options
# include the type of window to
# use and the number of points
# at which to estimate the density.
# n = the number of equally spaced
# points at which the density
# is evaluated.
```

1163

```
b1.boot <- storm.boot1$rep[,1]
b2.boot <- storm.boot1$rep[,2]

library(MASS)

par(fin=c(7.0,7.0),mex=1.5)

mm <- range(b1.boot)
min.int <- floor(mm[1])
max.int <- ceiling(mm[2])

truehist(b1.boot,xlim=c(min.int,max.int),
          density=.0001)

width.SJ(b1.boot)

[1] 0.6918326

b1.boot.dns1 <- density(b1.boot,n=200,width=.6)
b1.boot.dns2 <-
  list( x = b1.boot.dns1$x,
        y = dnorm(b1.boot.dns1$x,mean(b1.boot),
                  sqrt(var(b1.boot))))
```

1164

```

# Draw the non-parametric density

lines(b1.boot.dns1,lty=3,lwd=2)

# Draw normal densities

lines(b1.boot.dns2,lty=1,lwd=2)

legend(27.,-0.30,c("Nonparametric",
  "Normal approx."),
  lty=c(3,1),bty="n",lwd=2)

# Display the distribution of the estimate
# of the second parameter

par(fin=c(7.0,7.0),mex=1.5)

mm <- range(b2.boot)
min.int <- floor(mm[1])
max.int <- ceiling(mm[2])

```

1165

```

truehist(b2.boot,xlim=c(min.int,max.int),
  density=.001)

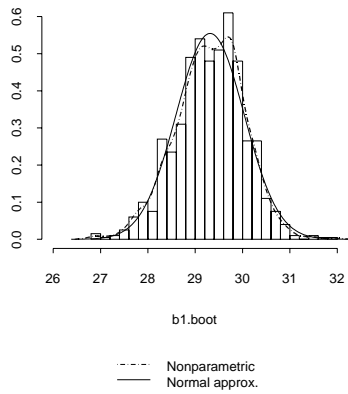
width.SJ(b2.boot)
[1] 0.7339957

b2.boot.dns1 <- density(b2.boot,n=200,width=.8)
b2.boot.dns2 <-
  list( x = b2.boot.dns1$x,
  y = dnorm(b2.boot.dns1$x,mean(b2.boot),
    sqrt(var(b2.boot))))

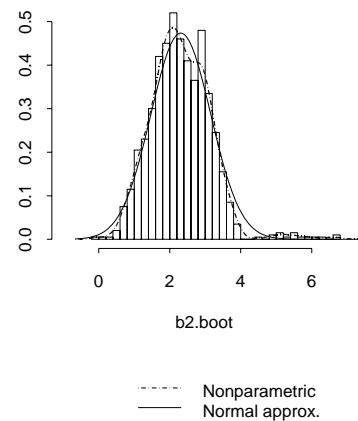
lines(b2.boot.dns1,lty=3,lwd=2)
lines(b2.boot.dns2,lty=1,lwd=2)
legend(1.,-0.30,c("Nonparametric",
  "Normal approx."),
  lty=c(3,1),bty="n",lwd=2)

```

1166



1167



1168

Bootstrap II: Fix the values of the explanatory variables

$$\{(W_j, V_j) : j = 1, \dots, n\}$$

- Compute residuals from fitting the model to the original sample

$$e_j = T_j - \frac{\bar{\beta}_1 V_j}{W_j - \bar{\beta}_2} \quad j = 1, \dots, n$$

- Approximate sampling from the population of random errors by taking a sample (with replacement) from $\{e_1, \dots, e_n\}$ say,

$$e_{1,b}^*, e_{2,b}^*, \dots, e_{n,b}^*$$

1169

- Create new observations:

$$(T_{j,b}^*, W_j, V_j) \quad j = 1, \dots, n$$

where

$$T_{j,b}^* = \frac{\bar{\beta}_1 V_j}{W_j - \bar{\beta}_2} + e_j^*$$

- Fit the model to the j -th bootstrap sample to obtain

$$\bar{\beta}_{1,b}^* \quad \text{and} \quad \bar{\beta}_{2,b}^*$$

- Repeat this for $b = 1, \dots, B$ bootstrap samples

1170

```
#####
# Bootstrap II:
#   Treat the regressors as fixed and
#   resample from the residuals
#####

# Center the residuals at zero and
# divide residuals by linear approximations
# to a multiple of the standard errors. These
# centered and scaled residuals approximately
# have the same first two moments as the
# random errors, but they are not quite
# uncorrelated.

rs <- scale(resid(storm.fm), center=T, scale=F)

grad.f <- deriv3( expr = Y ~ (b1*X1)/(X2-b2),
                  namevec = c("b1", "b2"),
                  function.arg = function(b1, b2, X1, X2, Y) NULL)
```

1171

```
g2 <- grad.f(b[1], b[2], stormer$Viscosity,
             stormer$Wt, stormer$Tim)
D <- attr(g2, "gradient")
h <- 1-diag(D%*%solve(t(D)%*%D)%*%t(D))
rs <- rs/sqrt(h)
dfe <- length(rs)-length(coef(storm.fm))
vr <- var(rs)
rs <- rs%*%sqrt(deviance(storm.fm)/dfe/vr)
```

1172

```

# Create a function to use in fitting
# the model to bootstrap samples

storm.bf2 <- function(rs)
  {assign("Tim", fitted(storm.fm) + rs,
    frame = 1)
  nls(formula = Tim ~ (b1*Viscosity)/(Wt-b2),
    data = stormer,
    start = coef(storm.fm)
  )$parameters }

summary(storm.fm)$parameters

      Value Std. Error  t value
b1 29.401328  0.9155353  32.113813
b2  2.218226  0.6655234   3.333054

```

1173

```

# Compute 1000 bootstrapped values of the
# regression parameters

```

```

storm.boot2 <- bootstrap(data=rs,
  statistic=storm.bf2(rs),B=1000)

```

```

Forming replications 1 to 100
Forming replications 101 to 200
Forming replications 201 to 300
Forming replications 301 to 400
Forming replications 401 to 500
Forming replications 501 to 600
Forming replications 601 to 700
Forming replications 701 to 800
Forming replications 801 to 900
Forming replications 901 to 1000

```

Call:

```

bootstrap(data = rs, statistic = storm.bf2(rs),
  B = 1000)

```

```

b1.boot <- storm.boot2$rep[,1]
b2.boot <- storm.boot2$rep[,2]

```

1174

```

# The BCA intervals may not be correctly
# computed by the following function

```

```

summary(storm.boot2)

```

Number of Replications: 1000

Summary Statistics:

	Observed	Bias	Mean	SE
b1	28.745	0.5719	29.317	0.8892
b2	2.432	-0.1682	2.264	0.6353

1175

Empirical Percentiles:

	2.5%	5%	95%	97.5%
b1	27.58	27.847	30.690	30.924
b2	1.03	1.241	3.295	3.459

BCa Confidence Limits:

	2.5%	5%	95%	97.5%
b1	26.236	26.267	29.698	29.985
b2	1.316	1.524	3.563	3.683

Correlation of Replicates:

	b1	b2
b1	1.0000	-0.9194

1176


```

# The following code is used to draw
# non-parametric estimated densities,
# normal densities, and histograms
# on the same graphic window.

par(fin=c(7.0,7.0),mex=1.3)

mm <- range(b1.boot)
min.int <- floor(mm[1])
max.int <- ceiling(mm[2])

truehist(b1.boot,xlim=c(min.int,max.int),
          density=.0001)

b1.boot.dns1 <- density(b1.boot,n=200,
                        width=width.SJ(b1.boot))
b1.boot.dns2 <-
  list( x = b1.boot.dns1$x,
        y = dnorm(b1.boot.dns1$x,mean(b1.boot),
                  sqrt(var(b1.boot))))

```

1177

```

# Draw a non-parametric density

lines(b1.boot.dns1,lty=3,lwd=2)

# draw normal densities

lines(b1.boot.dns2,lty=1,lwd=2)

legend(27,-0.22,c("Nonparametric",
                  "Normal approx."),
       lty=c(3,1),bty="n",lwd=2)

```

1178

```

# Display the distribution for the other
# parameter estimate

par(fin=c(7.0,7.0),mex=1.3)

mm <- range(b2.boot)
min.int <- floor(mm[1])
max.int <- ceiling(mm[2])

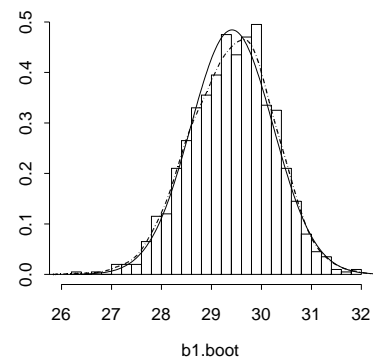
truehist(b2.boot,xlim=c(min.int,max.int),
          density=.00001)

b2.boot.dns1 <- density(b2.boot,n=200,
                        width=width.SJ(b2.boot))
b2.boot.dns2 <-
  list( x = b2.boot.dns1$x,
        y = dnorm(b2.boot.dns1$x,mean(b2.boot),
                  sqrt(var(b2.boot))))

lines(b2.boot.dns1,lty=3,lwd=2)
lines(b2.boot.dns2,lty=1,lwd=2)
legend(0.5,-0.25,c("Nonparametric",
                  "Normal approx."),lty=c(3,1),
       bty="n",lwd=2)

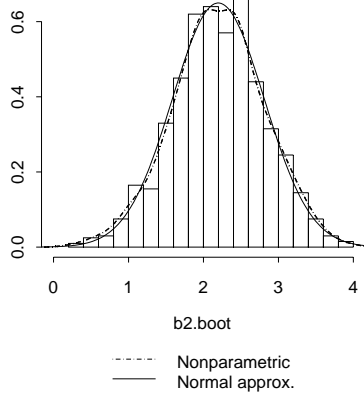
```

1179



----- Nonparametric
 _____ Normal approx.

1180



1181

Comparison of standard errors:

Para- meter	Asymptotic normal approx.	Random Bootstrap	Fixed Bootstrap
β_1	0.916	0.719	0.889
β_2	0.666	0.842	0.635

**Comparison of approximate 95%
confidence intervals:**

Para- meter	Asymptotic normal approx.	Random Bootstrap	Fixed Bootstrap
β_1	(27.50, 31.31)	(28.13, 30.40)	(27.42, 31.03)
β_2	(0.83, 3.60)	(1.06, 3.51)	(1.03, 3.46)

↑
 $\bar{\beta}_i \pm t_{21, .025} S_{\bar{\beta}_i}$

1182