

Section § 1.5 Bootstrap

8)

20) The bootstrap generates pseudodata

T_1^*, \dots, T_B^* for a statistic T_n that estimates θ ,

such that, under regularity conditions, applying certain large sample $100(1-\delta)\%$ PIs to the bootstrap sample results in a large sample $100(1-\delta)\%$ confidence intervals (CIs). The "*" means from the bootstrap resampling.

21) ^{PI9} Applying the shortest PI to T_1^*, \dots, T_B^* gives the shortest CI

22) ^{PI9} The "large sample" $100(1-\delta)\%$ percentile CI covers $\approx \lceil B(1-\delta) \rceil$ of the T_i^* .

Let $k_1 = \lceil B \delta/2 \rceil$ and $k_2 = \lceil B(1-\delta/2) \rceil$.

A common choice is $\left[T_{(k_1)}^*, T_{(k_2)}^* \right]$, which

is the nonparametric PI of [19] applied to the bootstrap sample. The shortest CI

is the "shortest" percentile CI that covers α cases.

23) Nonparametric bootstrap: Let the data

Y_1, \dots, Y_n be iid and let statistic

$T_n = t(Y_1, \dots, Y_n)$. Draw a sample of size n

Y_1^*, \dots, Y_n^* with replacement from Y_1, \dots, Y_n and

compute $T_1^* = t(Y_1^*, \dots, Y_n^*)$. Repeat B times

T_1^* generate \dots T_B^* bootstrap sample

24) know for E1 } For tiny B, given B bootstrap samples, compute

T_1^*, \dots, T_B^* and the bagging estimator

$$\bar{T}^* = \frac{1}{B} \sum_{i=1}^B T_i^* = \text{sample mean of the bootstrap samples}$$

Often $T_n = \bar{X}_n$, MED(n), range = $X_{(n)} - X_{(1)}$

ex) $Y_1, Y_2 = (0, 4)$ $T_n = \text{range} = X_{(n)} - X_{(1)} = 4 - 0 = 4$

bootstrap sample Y_i^*	T_i^*
0, 0	0
0, 4	4
4, 0	4
4, 4	0

$$\bar{T}^* = \frac{8}{4} = 2$$

see Example 2.10.

ex) data 1, 2, 5, 10, 50 MED(n) = 5 = T_n

bootstrap sample	ordered	T_i^*
2, 10, 1, 2, 2	1, 2, 2, 2, 10	$T_1^* = 2$
50, 10, 50, 2, 2	2, 2, 10, 50, 50	$T_2^* = 10$
10, 50, 2, 1, 1	1, 1, 2, 10, 50	$T_3^* = 2$

$$\bar{T}^* = \frac{\sum T_i^*}{B} = \frac{14}{3} = 4.6667$$

is $\left[\bar{T}^* - a_{(UB)}, \bar{T}^* + a_{(UB)} \right]$ is the closed interval centered at \bar{T}^* just large enough to cover $UB = \sqrt{B} (1.5)$ of the T_i^* where $a_i = |T_i^* - \bar{T}^*|$.

Replace T^* by T_n for the Bickel and Ren CI.

26) The percentile CI, CI in 25) and possibly the shorth CI are useful for robust statistics with good large sample theory and good bootstrap theory. $T_n = MED(n)$ is such a statistic.

27) often apply the bootstrap CIs for T_n and hope they simulate well. Then the CI and test are "ad hoc."

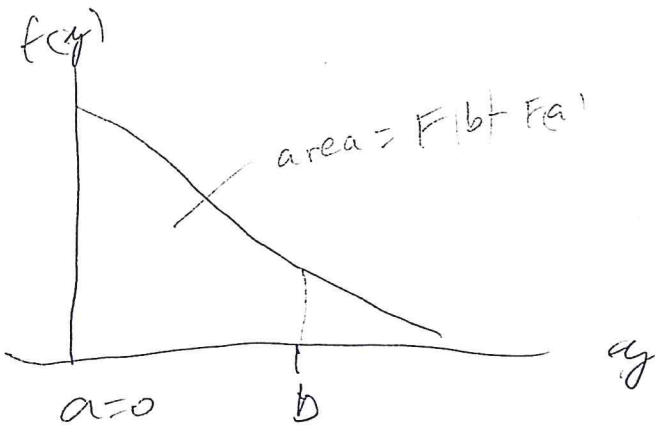
28) Bootstrap test: T_n estimates θ
 $\sqrt{n} (T_n - \theta) \xrightarrow{D} N(0, \sigma_T^2)$.

$H_0: \theta = \theta_0$ vs $H_A: \theta \neq \theta_0$

reject H_0 if θ_0 is not in the 100(1- α)% CI for θ .

distributions will be useful see p. 117, section

9.5



pdf of Y_T cdf $F(y|a,b)$ needs to integrate to 1.

Let $F(y) = P[Y \leq y]$. Use Assume $F(a) = F(a-)$

The truncated RV $Y_T = Y_T(a,b)$ with truncation points a and b has cdf

$$F_{Y_T}(y|a,b) = G(y) = \begin{cases} 0 & y < a \\ \frac{F(y) - F(a)}{F(b) - F(a)} & a \leq y \leq b \\ 1 & y > b \end{cases}$$

with mean $\mu_T = \mu_T(a,b) = \frac{\int_a^b y f(y) dy}{F(b) - F(a)}$

and variance $\sigma_T^2 = \sigma_T^2(a,b) = \frac{\int_a^b y^2 f(y) dy}{F(b) - F(a)} - \mu_T^2$

30} p. 21 The α trimmed mean

$$\bar{T}_n = \bar{T}_n(L_n, U_n) = \frac{1}{U_n - L_n} \sum_{i=L_n+1}^{U_n} Y(i) \quad \text{where}$$

$L_n = \lfloor n\alpha \rfloor$ and $U_n = n - L_n$. If $\alpha = 0.25$, say, the α trimmed mean is called the 25% trimmed mean.

A 95% CI for θ is $\bar{\theta} \pm t_{p, 0.975} SE(\bar{\theta})$

(10)

where $p = \text{degrees of freedom}$ is an integer ≥ 1 ,
 $P(t_p \leq t_{p, 0.975}) = 0.975$ and $t_{p, 0.975} \approx 1.96$

For $p > 30$,

12) Know for EI find $\hat{\theta}, p, SE(\hat{\theta})$ for 3 estimators

$\hat{\theta}$	p	$SE(\hat{\theta})$
I) \bar{Y}	$n-1$	s/\sqrt{n}
II) $M \pm D(n)$	$(n-L_n-1)$	$0.5[\bar{Y}_{(U_n)} - \bar{Y}_{(L_{n+1})}]$

where $U_n = n - L_n$ and $L_n = \lfloor \frac{n}{2} \rfloor - \lceil \frac{n}{4} \rceil$

Here $\lfloor x \rfloor = \text{greatest integer function}$ so $\lfloor 7.7 \rfloor = 7$

and $\lceil x \rceil = \text{smallest integer } \geq x$ so $\lceil 7.7 \rceil = 8$

$\hat{\theta}$	p	$SE(\hat{\theta})$
III) T_n (trimmed mean)	$U_n - L_n - 1$	$SE(T_n)$

13) If T_n is the 25% trimmed mean,

$$L_n = \lfloor \frac{n}{4} \rfloor, \quad U_n = n - L_n$$

$$T_n = \frac{1}{U_n - L_n} \sum_{i=L_n+1}^{U_n} Y(i) \text{ estimates } \theta = \mu_T = \int_{\gamma_{0.25}}^{\gamma_{0.75}} z \cdot g \text{ (rightly)}$$

Order data, delete smallest and largest 25%

take sample mean of remaining half set of data

$$\begin{cases} Y_{(i)} & L_{n+1} \leq i \leq U_n \\ Y_{(U_n)} & i \geq U_{n+1} \end{cases}$$

Let $S_n^2(d_1, \dots, d_n) =$ sample variance of the d_i

$$= \frac{\sum (d_i - \bar{d})^2}{n-1} = \frac{\sum d_i^2 - n(\bar{d})^2}{n-1}$$

Let $V_{sw}(L_n, U_n) = \frac{S_n^2(d_1, \dots, d_n)}{\left(\frac{U_n - L_n}{n}\right)^2}$

Then $SE(\bar{T}_n) = \sqrt{\frac{V_{sw}(L_n, U_n)}{n}}$

ex) data 66, 3, 8, 5, 2

sort data 2, 3, 5, 8, 66

earlier showed $\bar{Y} = 16.8$ $S = 27.5989$ $P = n-1 = 4$

so $SE(\bar{Y}) = \frac{S}{\sqrt{n}} = \frac{27.5989}{\sqrt{5}} = 12.343$

$MED(n) = 5$

$$L_n = \left\lfloor \frac{n}{2} \right\rfloor - \left\lceil \sqrt{\frac{n}{4}} \right\rceil = \left\lfloor \frac{5}{2} \right\rfloor - \left\lceil \sqrt{\frac{5}{4}} \right\rceil = \lfloor 2.5 \rfloor - \lceil 1.118 \rceil$$

$$= 2 - 2 = 0$$

$$\text{so } U_n = n - L_n = n = 5$$

$$SE(MED(n)) = 0.5 (Y_{(U_n)} - Y_{(L_{n+1})}) = \frac{Y_{(5)} - Y_{(1)}}{2}$$

$$= \frac{66 - 2}{2} = 32$$

$$P = U_n - L_n - 1 = 5 - 0 - 1 = 4$$

CI's based on $MED(n)$ tend to be too long for small n , works better for $n > 100$

$$L_n = \lfloor \frac{n+1}{2} \rfloor = \lfloor \frac{5}{2} \rfloor = 2, \quad U_n = n - L_n = 5 - 2 = 3$$

$$2, | 3, 5, 8, | 66$$

$$\bar{T}_n = \frac{3+5+8}{3} = \frac{16}{3} \approx \boxed{5.333}$$

$$d_i's \quad 3, 3, 5, 8, 8$$

$$\bar{d} = \frac{3+3+5+8+8}{5} = \frac{27}{5} = 5.4$$

$$S_n^2(d_i, \dots, d_n) = \frac{\sum d_i^2 - n(\bar{d})^2}{n-1} = \frac{171 - 5(5.4)^2}{4}$$

$$= \frac{25.2}{4} = 6.3$$

$$V_{sw} = \frac{6.3}{\left(\frac{4-1}{5}\right)^2} = 17.5$$

$$SE(\bar{T}_n) = \sqrt{\frac{V_{sw}}{n}} = \sqrt{\frac{17.5}{3}} = \sqrt{5.833} = \boxed{2.415}$$

$$p = U_n - L_n - 1 = 3 - 2 - 1 = \boxed{0}$$

using rpack

$$y \in C(66, 3, 8, 5, 2)$$

a) cci(y)

$$int = CI = (-17.469, 5.069)$$

← way too long

$$mean = \bar{y} = 16.8$$

$$SE = SE(\bar{y}) = 12.343$$

111 - 11 - [-83.846, 93.846] ← long 115
 MED = MED(n) = 5

SE_{bg} = SE(MED(n)) = 32

c) $\hat{\mu}_{CI}(y)$ CI = [-2.716, 13.383] good

check $T_n \pm t_{2, .975} SE(T_n) = \frac{16}{3} \pm 4.303(1.871)$

$\approx 5.3333 \pm 8.0509 = [-2.718, 13.384]$
 ↑
 + table, computer is more accurate

34) These are large sample CI's (n should be large, roughly want $p > 20$).

Assumptions i) Y_1, \dots, Y_n are iid with a pdf

ii) $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{D} N(0, \lim_{n \rightarrow \infty} n \text{var}(\hat{\theta}) = V_{\hat{\theta}})$

so $\hat{\theta} \approx N(\theta, \frac{V_{\hat{\theta}}}{n})$

often $\lim_{n \rightarrow \infty} n \text{var}(\hat{\theta}) = V_{\hat{\theta}}$
 eg $\text{var} \bar{Y} = \frac{\sigma^2}{n}$
 $SE(\bar{Y}) = \frac{\sigma}{\sqrt{n}}$

and $P(-z_{.975} < \frac{\hat{\theta} - \theta}{SE(\hat{\theta})} < z_{.975}) \approx .95$
 ↑
 for large n

or $P(-z_{.975} SE(\hat{\theta}) < \hat{\theta} - \theta < z_{.975} SE(\hat{\theta})) \approx .95$

or $P(-\hat{\theta} - z_{.975} SE(\hat{\theta}) < -\theta < -\hat{\theta} + z_{.975} SE(\hat{\theta})) \approx .95$

or $P(\hat{\theta} - z_{.975} SE(\hat{\theta}) < \theta < \hat{\theta} + z_{.975} SE(\hat{\theta})) \approx .95$

so $P\left(\hat{\theta} - z_{.975} SE(\hat{\theta}) < \theta < \hat{\theta} + z_{.975} SE(\hat{\theta})\right) \approx .95$

) For smaller n replacing $z_{.975}$ by $t_{p, .975}$ works ok so

$\hat{\theta} \pm t_{p, .975} SE(\hat{\theta})$ is a good CI.

ex) $\hat{\theta} = \bar{y}$ then $\sqrt{\text{var}(\bar{y})} = \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$

and $SE(\bar{y}) = \frac{s}{\sqrt{n}}$, $s^2 \approx \sigma^2$ (for large n)

(2.14.35) Simulation relies on computer generated pseudo random variables that can, for many purposes, be regarded as iid random variables Y_1, \dots, Y_n from a specified distribution, the "runs" are independent.

$Y_1, Y_2, Y_3 \sim N(0,1)$

ex) $rnorm(3)$ run 1: .861	.194	1.699	MED(Y_i)
$rnorm(3)$ run 2: -1.418	.579	.409	-.861
			-.409

ex) length of pregnancy $\approx N(\mu = 266, \sigma^2 = (16)^2)$

$rnorm(3, mean=266, SD=16)$

run 1	266.34	281.46	251.898	max 281.46
run 2	261.81	232.64	266.07	

start a year, what will be the longest pregnancy? (12.5)

To simulate 100 such hospitals use the command

$X \leftarrow \text{matrix}(\text{rnorm}(100000, \text{mean}=266, \text{sd}=16), \text{nrow}=100, \text{ncol}=1000)$

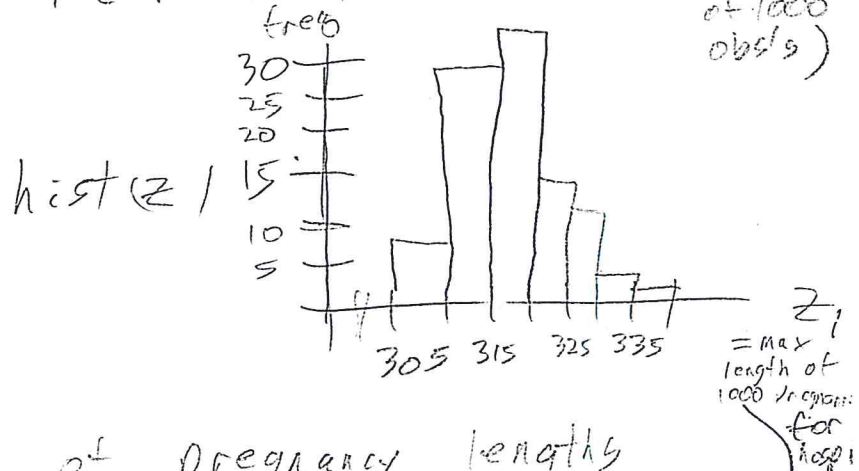
X is 100×1000 and each row represents a hospital with 1000 pregnancy lengths.

$Z \leftarrow \text{apply}(X, 1, \text{max})$

$\text{hist}(X)$ to see all births at the 100 hospitals

find the max of each row (= sample of 1000 obs's)

min(Z) 307.76
max(Z) 337.83
mean(Z) 317.70
median(Z) 316.52
sqrt(var(Z)) 6.39
mad(Z, constant=1) 4.15



(Actually the right tail of pregnancy lengths may not be normal. - Hospital might induce labour or do a csection for 30 or 50 days overdue.)

26) The point of simulation is that you can do large numbers of computer experiments often in seconds.

37) Applications: i) check computer programs: does simulation match theory
ii) demonstrate 10 or 20 graphs in rapid succession so user can see typical graph

ex HW 1.3d 4 histograms of normal data with sample size: $n=100$. Most of the histograms do