

1.5.2 Basic Bootstrap CI (Corrects for bias).

based on residuals.

The $100(1 - \alpha)\%$ Basic Bootstrap CI for θ is

$$\left(\hat{\theta} - [\hat{\theta}_{1-\alpha/2} - \hat{\theta}], \hat{\theta} - [\hat{\theta}_{\alpha/2} - \hat{\theta}] \right) = \left(2\hat{\theta} - \hat{\theta}_{1-\alpha/2}, 2\hat{\theta} - \hat{\theta}_{\alpha/2} \right).$$

↑ estimate from original sample
↑ $1-\alpha/2$ quantile based on bootstrap distribution

Why?

Let $\varepsilon = \hat{\theta} - \theta$ and $\varepsilon_{1-\alpha/2}, \varepsilon_{\alpha/2}$ quantiles of dsn of ε 's.

$$P(\varepsilon_{\alpha/2} \leq \hat{\theta} - \theta \leq \varepsilon_{1-\alpha/2}) = 1 - \alpha \Rightarrow (1 - \alpha)\% \text{ CI is } \left(\hat{\theta} - \varepsilon_{1-\alpha/2}, \hat{\theta} - \varepsilon_{\alpha/2} \right).$$

Assumptions/usage

$$\varepsilon_{\alpha/2} \approx \hat{\theta}_{\alpha/2} - \hat{\theta} \text{ bootstrap version.}$$

- Corrects for bias, but slightly harder to explain than percentile.
- Not transformation invariant.

1.5.3 Bootstrap t CI (Studentized Bootstrap) Consider $Z = \frac{\hat{\theta} - E(\hat{\theta})}{se(\hat{\theta})}$

Even if the distribution of $\hat{\theta}$ is Normal and $\hat{\theta}$ is unbiased for θ , the Normal distribution is not exactly correct for z .

because we need to estimate $se(\hat{\theta})$. $\Rightarrow t^* = \frac{\hat{\theta} - E(\hat{\theta})}{\hat{se}(\hat{\theta})} \sim t_{n-1}$? No

\uparrow
where $\hat{se}(\hat{\theta}) = sd(\hat{\theta}^{*(1)}, \dots, \hat{\theta}^{*(B)})$

Additionally, the distribution of $\hat{se}(\hat{\theta})$ is unknown.

So we cannot just claim $t^* \sim t_{n-1}$.

\Rightarrow The bootstrap t interval does not use a Student t distribution as the reference distribution, instead we estimate the distribution of a "t type" statistic by resampling.

The $100(1 - \alpha)\%$ Bootstrap t CI is based on bootstrap samples $\hat{\theta}^{*(1)}, \dots, \hat{\theta}^{*(B)}$

$(\hat{\theta} - t_{1-\alpha/2}^* \hat{se}(\hat{\theta}), \hat{\theta} + t_{\alpha/2}^* \hat{se}(\hat{\theta}))$

estimate from original data sample (pointing to $\hat{\theta}$)

quantile of the bootstrap dist of our "t-type" statistic. (pointing to $t_{1-\alpha/2}^*$ and $t_{\alpha/2}^*$)

Overview

t-type statistic: $t^{(1)} = \frac{\hat{\theta}^{*(1)} - \hat{\theta}}{\hat{se}(\hat{\theta}^{*(1)})}, \dots, t^{(B)} = \frac{\hat{\theta}^{*(B)} - \hat{\theta}}{\hat{se}(\hat{\theta}^{*(B)})}$

= bootstrap estimate of $se(\hat{\theta}^{(1)})$ based on the first bootstrap sample?? Double bootstrap!*

To estimate the "t style distribution" for θ ,

1. Compute $\hat{\theta}$
2. For $b=1, \dots, B$
 - (a) Sample w/ replacement
 $y^{(b)} = (y_1^{(b)}, \dots, y_n^{(b)})$
 - (b) compute $\hat{\theta}^{*(b)}$
 - (c) For each $r=1, \dots, R$
 - (i) Sample w/ replacement from $y^{(b)}$
 $y^{(b)(r)} = (y_1^{(b)(r)}, \dots, y_n^{(b)(r)})$
 - (ii) compute $\hat{\theta}^{*(b)(r)}$
 - (d) compute $\hat{se}(\hat{\theta}^{*(b)}) = sd(\hat{\theta}^{*(b)(1)}, \dots, \hat{\theta}^{*(b)(R)})$
 - (e) compute $t^{(b)} = \frac{\hat{\theta}^{*(b)} - \hat{\theta}}{\hat{se}(\hat{\theta}^{*(b)})}$
3. get quantiles $t_{\alpha/2}^*, t_{1-\alpha/2}^*$
4. compute CI.

Assumptions/usage

- Computationally intensive.
- = Not doing anything for bias/skewness.
- Need $\hat{\theta}$ independent of $\hat{SE}(\hat{\theta})$.

This idea is based on "pivot quantities" = a function of observations and parameters whose probability distn doesn't depend on the parameter.

$$\text{E.g., } \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$$

↑ no μ in here!

This can help us to obtain a bootstrap CI for μ .

You could create other pivot quantities to extend via bootstrap, e.g.

$$\text{If } Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2) \quad \frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{n-1}$$

$$\Rightarrow \left(\frac{(n-1)S^2}{\chi^2_{1-\alpha/2}}, \frac{(n-1)S^2}{\chi^2_{\alpha/2}} \right) \text{ would be a } 95\% \text{ CI for } \sigma^2$$

Bootstrap version does not assume $Y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$. Now need to estimate χ^2 -type quantile using the bootstrap.

1.5.4 BCa CIs

"bias-corrected accelerated"

correct for skew.

Modified version of percentile intervals that adjusts for bias of estimator and skewness of the sampling distribution.

This method automatically selects a transformation so that the normality assumption holds.

Idea:

Assume there exists a monotonically increasing function g and constants a and b s.t.

$$U = \frac{g(\hat{\theta}) - g(\theta)}{1 + ag(\theta)} + b \sim N(0,1) \text{ where } 1 + ag(\theta) > 0.$$

By the bootstrap principle

$$U^* = \frac{g(\hat{\theta}^*) - g(\hat{\theta})}{1 + ag(\hat{\theta})} + b \overset{\text{approx}}{\sim} N(0,1).$$

⇒ For any quantile of a standard Normal dsn,

$$\begin{aligned} \alpha &\approx P^* [U^* \leq z_\alpha] \\ &= P^* \left[\frac{g(\hat{\theta}^*) - g(\hat{\theta})}{1 + ag(\hat{\theta})} + b \leq z_\alpha \right] \\ &= P^* \left[\hat{\theta}^* \leq \underbrace{g^{-1} \left(g(\hat{\theta}) + (z_\alpha - b)(1 + ag(\hat{\theta})) \right)}_{\hat{\theta}_\alpha} \right] \end{aligned}$$

The α quantile from the bootstrap dsn of $\hat{\theta}^*$, denoted $\hat{\theta}_\alpha$, is observable from BS dsn.

$$\Rightarrow g^{-1} \left(g(\hat{\theta}) + (z_\alpha - b)(1 + ag(\hat{\theta})) \right) \approx \hat{\theta}_\alpha$$

To use this, consider U : $1 - \alpha = P(U > z_\alpha)$

$$= P \left(\theta < \underbrace{g^{-1} \left(g(\hat{\theta}) + \frac{b - z_\alpha}{1 - a(b - z_\alpha)} [1 + ag(\hat{\theta})] \right)} \right)$$

Notice similarity to above.

⇒ If we could find β such that $\frac{b - z_\alpha}{1 - a(b - z_\alpha)} = z_\beta - b$ then the bootstrap principle can be applied to conclude $\theta < \hat{\theta}_\beta$ will be an appropriate $1 - \alpha$ upper CI limit.

$$\Rightarrow z_\beta = \frac{b - z_\alpha}{1 - a(b - z_\alpha)} + b \Rightarrow \beta = \Phi \left(\frac{b + z_{1-\alpha}}{1 - a(b + z_\alpha)} + b \right).$$

2-sided argument is similar.

⇒ β^{th} quantile of $N(0,1)$ ⇒ ① find a, b , ② compute β , ③ Find β^{th} quantile of empirical dsn of $\hat{\theta}^*$

b : let p_0 denote the fraction of obs. from bootstrap dsn s.t. $\hat{\theta}^* < \hat{\theta}$. Since g is monotone, this is the same fraction of $\hat{\theta}^*$'s s.t. $g(\hat{\theta}^*) < g(\hat{\theta})$.

⇒ $P(Z < b) = p_0$ where $Z \sim N(0,1)$ gives us a way to estimate b (If the bootstrap dsn has $\hat{\theta}$ as its median, then $b=0$).

⇒ b corrects for bias.

a : let $S_{-i} = \{y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n\}$ and let $\hat{\theta}_{-i}$ denote the estimate of θ based on S_{-i} : $a = \frac{\sum_{i=1}^n (\hat{\theta}_{(i)} - \hat{\theta}_{-i})^3}{6 \left[\sum_{i=1}^n (\hat{\theta}_{(i)} - \hat{\theta}_{-i})^2 \right]^{3/2}}$ ← based on skewness estimator where $\hat{\theta}_{(i)}$ is the mean of $\hat{\theta}_i$'s.

If $\alpha=0$, don't adjust for skewness \Rightarrow BC interval.

17

The BCa method uses bootstrapping to estimate the bias and skewness then modifies which percentiles are chosen to get the appropriate confidence limits for a given data set.

In summary,

BCa is like the percentile bootstrap CI, but instead of $(\hat{\theta}_{0.1/2}, \hat{\theta}_{1-0.1/2})$, choose better quantiles to account for bias and skewness.

Has better coverage than percentile method in empirical studies, but harder to explain.

Your Turn

We will consider a telephone repair example from Hesterberg (2014). `verizon` has repair times, with two groups, CLEC and ILEC, customers of the “Competitive” and “Incumbent” local exchange carrier.

```
library(resample) # package containing the data
```

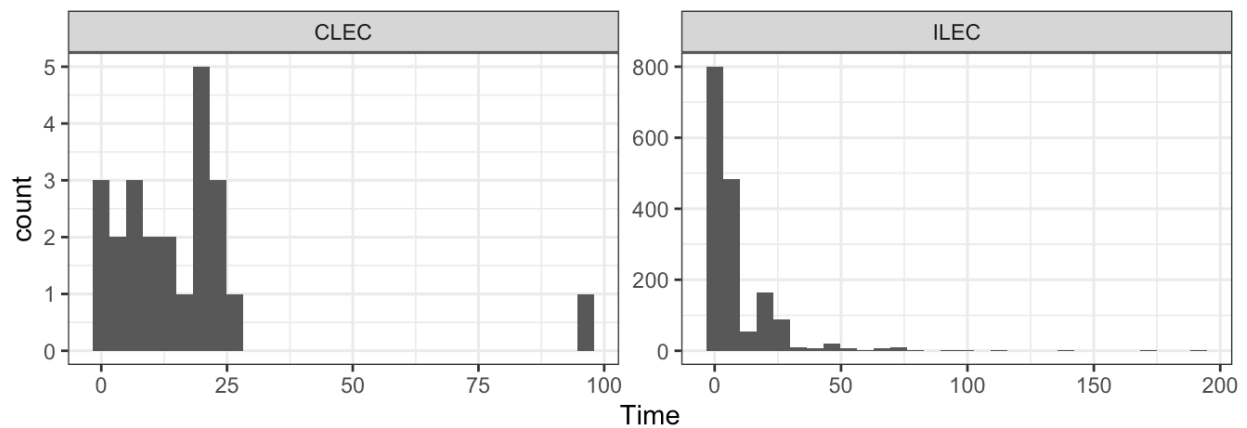
```
data(Verizon)
head(Verizon)
```

```
##      Time Group
## 1 17.50  ILEC
## 2  2.40  ILEC
## 3  0.00  ILEC
## 4  0.65  ILEC
## 5 22.23  ILEC
## 6  1.20  ILEC
```

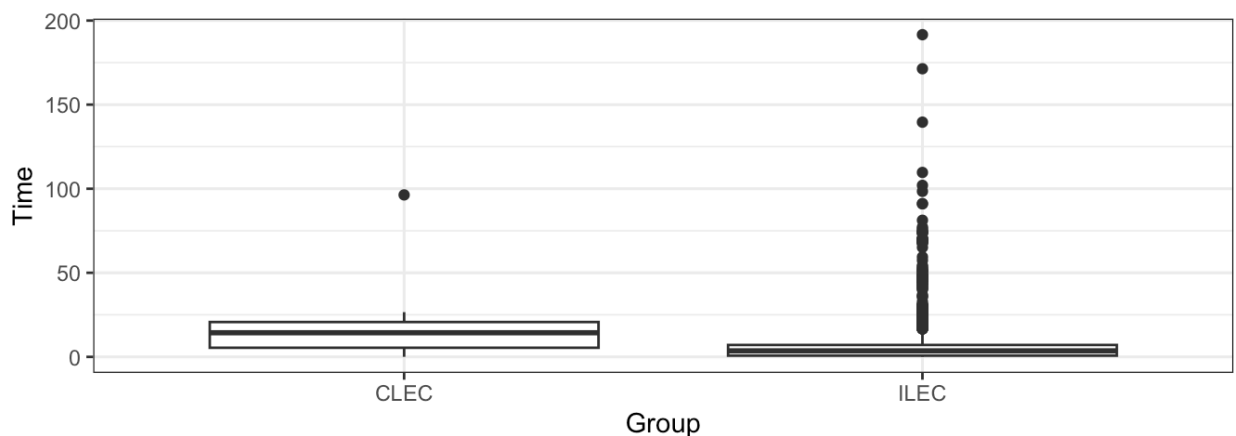
```
Verizon |>
  group_by(Group) |>
  summarize(mean = mean(Time), sd = sd(Time), min = min(Time), max =
             max(Time)) |>
  kable()
```

Group	mean	sd	min	max	
CLEC	16.509130	19.50358	0	96.32	23
ILEC	8.411611	14.69004	0	191.60	1664

```
ggplot(Verizon) +
  geom_histogram(aes(Time)) +
  facet_wrap(~Group, scales = "free")
```



```
ggplot(Verizon) +
  geom_boxplot(aes(Group, Time))
```



1.6 Bootstrapping CIs

→ you could write your own.

There are many bootstrapping packages in R, we will use the boot package. The function `boot` generates R resamples of the data and computes the desired statistic(s) for each sample. This function requires 3 arguments:

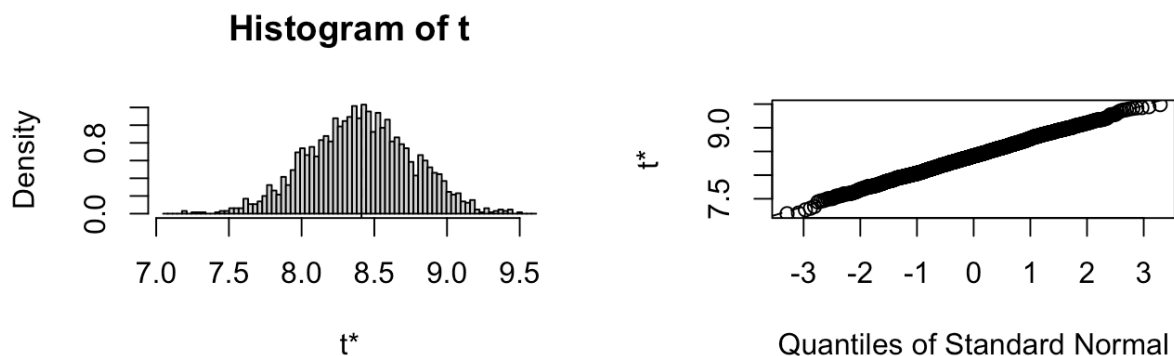
1. $data$ = the data from the original sample (data.frame or matrix).
2. $statistic$ = a function to compute the statistic from the data where the first argument is the data and the second argument is the indices of the observations in the bootstrap sample.
3. R = the number of bootstrap replicates.

```
library(boot) # package containing the bootstrap function

mean_func <- function(x, idx) {
  mean(x[idx])
}

ilec_times <- Verizon[Verizon$Group == "ILEC",]$Time
boot.ilec <- boot(ilec_times, mean_func, 2000)

plot(boot.ilec)
```



If we want to get Bootstrap CIs, we can use the `boot.ci` function to generate the different nonparametric bootstrap confidence intervals.



```
boot.ci(boot.ilec, conf = .95, type = c("perc", "basic", "bca"))
```

```
## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 2000 bootstrap replicates
##
## CALL :
## boot.ci(boot.out = boot.ilec, conf = 0.95, type = c("perc",
## "basic",
## "bca"))
##
## Intervals :
## Level      Basic          Percentile          BCa
## 95%    ( 7.733, 9.110 )  ( 7.714, 9.091 )  ( 7.755, 9.125 )
## Calculations and Intervals on Original Scale
```



```
## we can do some of these on our own
## percentile
quantile(boot.ilec$t, c(.025, .975))
```

```
##      2.5%      97.5%
## 7.714075 9.084725
```

```
## basic
2*mean(ilec_times) - quantile(boot.ilec$t, c(.975, .025))
```

```
##      97.5%      2.5%
## 7.738496 9.109147
```

To get the studentized bootstrap CI, we need our statistic function to also return the variance of $\hat{\theta}$.

```
mean_var_func <- function(x, idx) {
  c(mean(x[idx]), var(x[idx])/length(idx))
}

boot.ilec_2 <- boot(ilec_times, mean_var_func, 2000)
boot.ci(boot.ilec_2, conf = .95, type = "stud")
```

```
## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 2000 bootstrap replicates
##
## CALL :
## boot.ci(boot.out = boot.ilec_2, conf = 0.95, type = "stud")
##
## Intervals :
## Level      Studentized
## 95%      ( 7.728,  9.183 )
## Calculations and Intervals on Original Scale
```

Which CI should we use?