

Bootstrap (Recap)

- $X_1, X_2, \dots, X_n \sim \text{iid } F \text{ (CDF)}$
- $\theta = \text{parameter of interest} - \text{feature of } F$
- $\hat{F} = \text{estimated CDF (parametric or nonparametric)}$
- $\hat{\theta} = \text{estimator of } \theta - \text{same feature of } \hat{F}$
- Bootstrap resample: $X_1^*, X_2^*, \dots, X_n^* \sim \text{iid } \hat{F}$

sample	population	Estimator
original	F	$\hat{\theta}$
Bootstrap	\hat{F}	$\hat{\theta}^*$

- Generate a large # (b) of resamples:

Resample #	Resample	Estimator	bootstrap dist. of $\hat{\theta}$
1	$X_{11}^*, X_{12}^*, \dots, X_{1n}^*$	$\hat{\theta}_1^*$	
2	$X_{21}^*, X_{22}^*, \dots, X_{2n}^*$	$\hat{\theta}_2^*$	
\vdots	\vdots	\vdots	
b	$X_{b1}^*, X_{b2}^*, \dots, X_{bn}^*$	$\hat{\theta}_b^*$	

- If n is small; take $b = \# \text{ possible resamples}$.

- To estimate a feature η of sampling dist. of $\hat{\theta}$, we $\eta^* = \text{same feature estimated from the bootstrap draws of } \hat{\theta}^*$

- Essentially: "Replace θ by $\hat{\theta}$, and $\hat{\theta}$ by $\hat{\theta}^*$ "

- If $\eta = \alpha\text{-th quantile of } \hat{\theta} - \theta$; then $\hat{\eta}^* = \alpha\text{-th sample percentile of } \hat{\theta}_1^* - \hat{\theta}, \dots, \hat{\theta}_b^* - \hat{\theta}$
 $= \hat{\theta}_{((b+1)\alpha)}^* - \hat{\theta}$

Four Bootstrap CIs for θ

1. Normal approximation CI: Use z critical point but

correct $\hat{\theta}$ for bias and use \hat{V}^* to estimate V .

- Estimated bias of $\hat{\theta} \notin \hat{\theta}^*$ $= \frac{1}{b} \sum_{k=1}^b \hat{\theta}_k^* - \hat{\theta}$
- CI: $[\hat{\theta} - \hat{B}^* - z_{1-\alpha/2} \widehat{SE}^*, \hat{\theta} - \hat{B}^* - z_{\alpha/2} \widehat{SE}^*]$.

2. Studentized bootstrap CI: Use bootstrap critical point of T instead of z critical point.

$$T = \frac{\hat{\theta} - \theta}{SE(\hat{\theta})}$$

- Get $T_1^* = (\hat{\theta}_1^* - \hat{\theta}) / \widehat{SE}_1^*, \dots, T_b^* = (\hat{\theta}_b^* - \hat{\theta}) / \widehat{SE}_b^*$
- Estimated α -th percentile of $\underline{T} = T_{((b+1)\alpha)}^*$
- CI: $[\hat{\theta} - t_{((b+1)(1-\alpha/2))}^* \widehat{SE}, \hat{\theta} - t_{((b+1)(\alpha/2))}^* \widehat{SE}]$.

3. Basic bootstrap CI: Based on percentiles of $\hat{\theta} - \theta$ rather than $(\hat{\theta} - \theta)/\widehat{SE}$. Use bootstrap to estimate them. Notice

$$1 - \alpha = P(a_{\alpha/2} \leq \hat{\theta} - \theta \leq a_{1-\alpha/2})$$

for percentile of $\hat{\theta} - \theta$

$$= P[a_{\alpha/2} - \hat{\theta} \leq -\theta \leq a_{1-\alpha/2} - \hat{\theta}]$$

$$= P[\cancel{a_{\alpha/2} - \hat{\theta}} \leq \hat{\theta} - \cancel{a_{1-\alpha/2}}] \leq \theta \leq \hat{\theta} - \cancel{a_{\alpha/2}}]$$

$$\bullet \text{ Estimated } a_{\alpha} = \hat{\theta}^*_{((b+1)\alpha)} - \hat{\theta} \rightarrow \hat{\theta}^*_{((b+1)(1-\alpha/2))} - \hat{\theta}$$

$$\bullet \text{ CI: } \left[\hat{2\theta} - \hat{\theta}^*_{((b+1)(1-\alpha/2))}, \hat{2\theta} - \hat{\theta}^*_{((b+1)(\alpha/2))} \right]$$

• Doesn't require SE .

4. Percentile bootstrap CI: Works as in basic bootstrap \rightarrow known f_n .
 but uses "magic." Suppose there exists a transformation h so
 that the distribution of $h(\hat{\theta}) - h(\theta)$ is symmetric about zero. [typically
 of f_n .]
 Let $U = h(\hat{\theta})$. As before, we can write

$$\begin{aligned}
 1 - \alpha &= P(a_{\alpha/2} \leq U - h(\theta) \leq a_{1-\alpha/2}) \\
 &\stackrel{\text{symmetry}}{=} P(-a_{1-\alpha/2} \leq U - h(\theta) \leq -a_{\alpha/2}) \\
 &= P(U + a_{\alpha/2} \leq h(\theta) \leq U + a_{1-\alpha/2})
 \end{aligned}$$

\uparrow exists for $h(\theta)$.

- Estimated $\underline{a}_\alpha = U_{((b+1)\alpha)}^* - U$
- $U + a_{\alpha/2} \approx U + \left\{ U_{((b+1)(\alpha/2))}^* - U \right\} = U^*_{((b+1)(\alpha/2))} - a_{1-\alpha/2}$
- Similarly, $U + a_{1-\alpha/2} = U^*_{((b+1)(1-\alpha/2))}$

Therefore,

$$\begin{aligned}
 1 - \alpha &\approx P\left(U^*_{((b+1)(\alpha/2))} \leq h(\theta) \leq U^*_{((b+1)(1-\alpha/2))}\right) \\
 &= P\left[h^{-1}\left\{U^*_{((b+1)(\alpha/2))}\right\} \leq h^{-1}h(\theta) \leq h^{-1}\left\{U^*_{((b+1)(1-\alpha/2))}\right\}\right] \\
 &= P\left[\theta^*_{((b+1)(\alpha/2))} \leq \theta \leq \theta^*_{((b+1)(1-\alpha/2))}\right]
 \end{aligned}$$

- CI: $\left[\hat{\theta}_{((b+1)(\alpha/2))}^*, \hat{\theta}_{((b+1)(1-\alpha/2))}^* \right]$.

- Magic: No need to know ψ' .

Q. Which method to use?

Research shows that studentized bootstrap is the best choice, but it requires \widehat{SE} . However, if \widehat{SE} is not available, then percentile bootstrap is often the next best choice. More accurate versions of this method are available.

Example: Recall the CPU time data from Example 8.12 on page 217. We had seen that a gamma distribution fit well to these data. Suppose we would like to perform inference on median cpu time.

R code:

```
# use install.packages("boot") to first install  
# the package and then load it
```

```
library(boot)
```

```
# read the cpu data (we have seen these before)
```

```
> (cpu <- scan(file="cputime.txt"))
```

```
Read 30 items
```

```
[1] 70 36 43 69 82 48 34 62 35 15 59 139  
46 37 42 30 55 56
```

```
[19] 36 82 38 89 54 25 35 24 22 9 56 19
>
```

```
# Parameter of interest: Median
```

```
#####
```

```
# Nonparametric Bootstrap #
```

```
#####
```

```
median.npar <- function(x, indices) {
  result <- median(x[indices])
  return(result)
}
```

```
> (median.npar.boot <- boot(cpu, median.npar, R=999,
sim="ordinary", stype="i"))
```

```
ORDINARY NONPARAMETRIC BOOTSTRAP
```

Call:

```
boot(data = cpu, statistic = median.npar, R = 999,  
sim = "ordinary", stype = "i")
```

Bootstrap Statistics :

	original	bias	std. error
t1*	42.5	0.6721722	5.876943

>

*both wanted
using
bootstrap.*

Let's verify the calculations

See what ~~is~~ else is stored in median.npar.boot

```
> names(median.npar.boot)
```

```
[1] "t0"
```

```
      "t"
```

```
      "data"
```

sample median

*bootstrap
median*


```
"seed"      "statistic"
[7] "sim"      "call"      "stype"      "strata"
"weights"
>
```

```
> median(cpu)
[1] 42.5
```

```
> median.npar.boot$t0
[1] 42.5
```

```
> mean(median.npar.boot$t)-median.npar.boot$t0
```

```
[1] 0.6721722 — bias } reported by "boot"
```

```
> sd(median.npar.boot$t)
```

```
[1] 5.876943 — se — No warning!
```

Warning message:

sd(<matrix>) is deprecated.

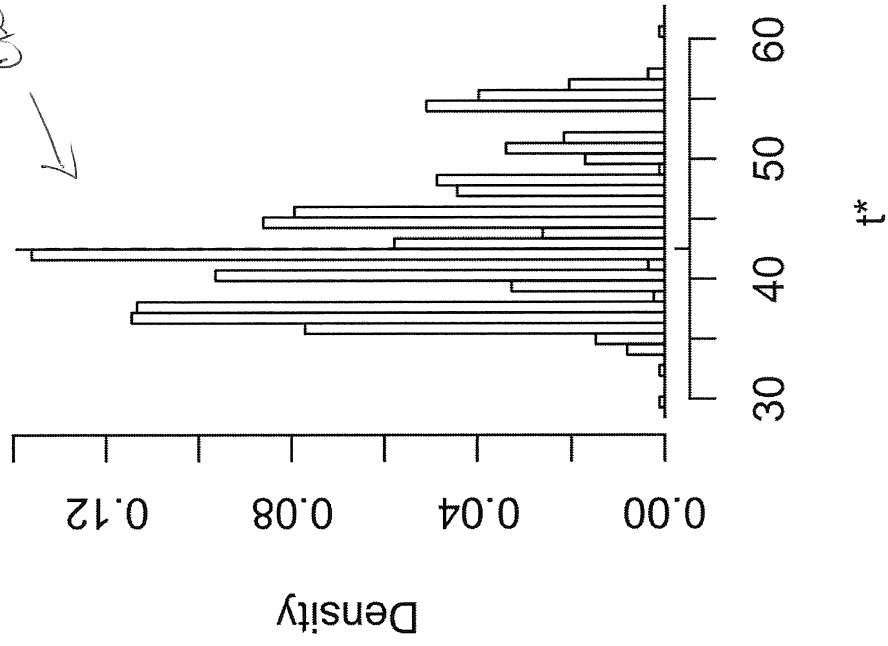
```
Use apply(*, 2, sd) instead.
```

```
>
```

```
# See the bootstrap distribution of median estimate
```

```
plot(median.npar.boot)
```

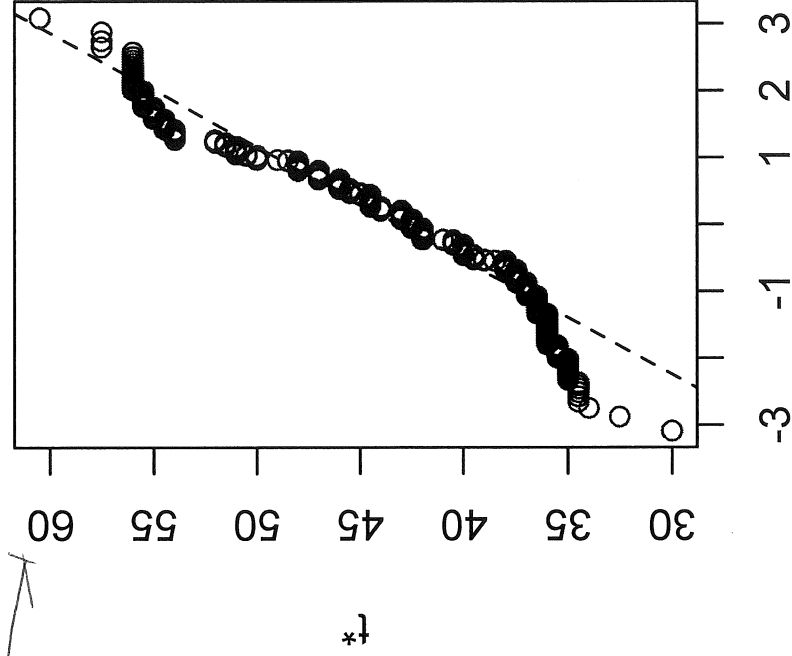
Histogram of t



clearly

not normal

⇒ the usual z-interval won't be accurate



Quantiles of Standard Normal

```
# Get the 95% confidence interval for median
```

```
> boot.ci(median.npar.boot)
```

```
BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
```

```
Based on 999 bootstrap replicates
```

```
CALL :
```

```
boot.ci(boot.out = median.npar.boot)
```

```
Intervals :
```

Level	Normal	Basic
95%	(30.31, 53.35)	(29.50, 49.50)

Level	Percentile
-------	------------

95%	(35.5, 55.5)	(35.0, 55.5)
-----	---------------	---------------

Calculations and Intervals on Original Scale

Warning message:

BCa → "not from this course"

```
In boot.ci(median.npar.boot) :  
bootstrap variances needed for studentized intervals  
>
```

```
# Let's verify
```

```
# Normal approximation method
```

```
> c(42.5 - 0.6721722 - qnorm(0.975) * 5.876943,  
    42.5 - 0.6721722 - qnorm(0.025) * 5.876943)  
[1] 30.30923 53.34642
```

```
>
```

```
# Percentile bootstrap method
```

```
> sort(median.npar.boot$t)[c(25, 975)]  
[1] 35.5 55.5
```

```
>
```

Everything
with
formulas.
matches

```
# Basic bootstrap method
```

```
> c(2*42.5-55.5, 2*42.5-35.5)
[1] 29.5 49.5
```

```
>
```

```
#####
```

```
# Parametric Bootstrap #
```

```
#####
```

```
# Saw earlier that a Gamma distribution fit well
# to these data
```

```
# > ml.est$par
```

```
# [1] 3.63149628 0.07529459
```

```
# >
```

```
# first element = shape parameter mle,
```

*parameter value -
"beta" is good*



```
# second element = rate parameter mle
```

$f = \text{gamma}(v, 1)$

```
median.par <- function(x) { return(median(x)) }
```

```
resamp.par <- function(x, mle) {
```

```
  xsim <- rgamma(length(x), mle$shape, mle$rate)
```

```
  return(xsim)
```

```
}
```

$\hat{\text{median}}^* = \text{MLE of median of gamma dist.}$
 $= \text{gamma}(0.50, 1)$
 $(\hat{\text{shape}}^*, \hat{\text{rate}}^*)$

```
> (median.par.boot <- boot(cpu, median.par, R=999,  
  sim="parametric", ran.gen=resamp.par,  
  mle=list(shape=3.63149628, rate=0.07529459)))
```

PARAMETRIC BOOTSTRAP

Call:

```
boot(data = cpu, statistic = median.par, R = 999, sim = "p  
ran.gen = resamp.par, mle = list(shape = 3.63149628, r
```

Bootstrap Statistics :

	original	bias	std. error
t1*	42.5	1.696431	5.220725

>

boot uses sample median as estimator rather than its
MLE which is the median of the fitted gamma distribution

```
> qgamma(0.5, shape = 3.63149628, rate = 0.07529459)
```

```
[1] 43.88304
```

>

↓ MLE of median
based on original sample


```
> boot.ci(boot.ci(median.par.boot))
Error in empinf(boot.out, index = index, t = t.o, ...) :
  influence values cannot be found from a parametric bootstrap
In addition: Warning message:
In boot.ci(median.par.boot) :
  bootstrap variances needed for studentized intervals
>
```

- boot package does not seem to work well for parametric bootstrap. But that's okay — this will be your Mini-
Project 4 extra credit assignment.