

Statistics for Data Science, Winter 2023

Chapter 3:
Bootstrap and Jackknife

Bootstrap and Jackknife are simulation based methods for inference/uncertainty quantification.

CIs and testing

pros:

- ▶ Conceptually simple
- ▶ No need for mathematically complex probability calculations

cons:

- ▶ Computationally expensive
- ▶ Requires large sample size (as for the normal approximation)
- ▶ Can go wrong even in large samples!
- ▶ Use with care!

RECALL: SAMPLING DISTRIBUTION



$X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta), \theta \in (0, 1), \hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$

$$\frac{\hat{\theta}_n - \theta}{se_\theta} \xrightarrow[n \rightarrow \infty]{d.} N(0, 1), \quad se_\theta = \sqrt{\text{Var}_\theta[\hat{\theta}_n]} = \sqrt{\frac{\theta(1-\theta)}{n}}$$

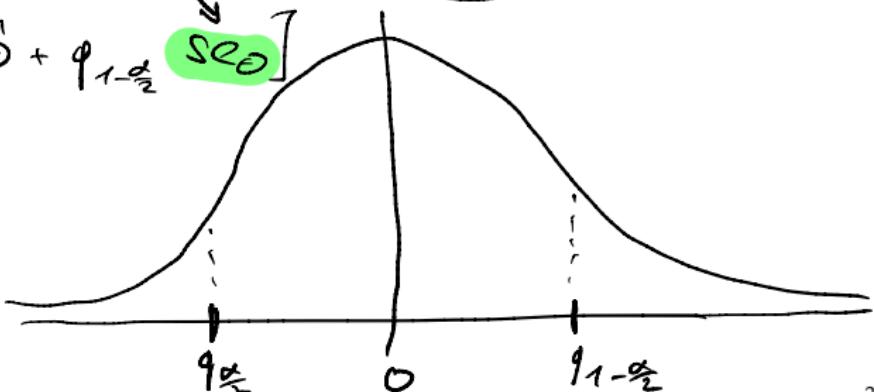
estimate se_θ by plug-in rule $se_{\hat{\theta}_n} = \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}}$

$$P\left(q_{\frac{\alpha}{2}} \leq \frac{\hat{\theta}_n - \theta}{se_{\hat{\theta}_n}} \leq q_{1-\frac{\alpha}{2}}\right) \approx 1-\alpha \sqrt{\frac{\hat{\theta}_n - \theta}{se_{\hat{\theta}_n}}} \approx N(0, 1)$$

\nwarrow unknown \searrow

$$[\hat{\theta}_n - q_{1-\frac{\alpha}{2}} se_{\hat{\theta}_n}, \hat{\theta}_n + q_{1-\frac{\alpha}{2}} se_{\hat{\theta}_n}]$$

$$\begin{aligned}\hat{se}_{\hat{\theta}_n} &= se_{\hat{\theta}_n} = \\ &= \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}}\end{aligned}$$



IN GENERAL



$$X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta, \theta \in \Theta,$$

We often have a CLT:

$$\frac{\hat{\theta}_n - \theta}{se_\theta} \xrightarrow[n \rightarrow \infty]{d.} N(0, 1)$$

But the standard error $se_\theta := \sqrt{\text{Var}_\theta[\hat{\theta}_n]}$ may be hard to compute or the plug-in idea may not be feasible.

First Bootstrap idea: *se bootstrap*

Estimate se_θ (without knowing its exact analytical form) and rely on the CLT.

Second Bootstrap idea: *pivotal bootstrap*

Estimate quantiles of the sampling distribution of $\hat{\theta}_n - \theta$. No Gaussian approximation required.

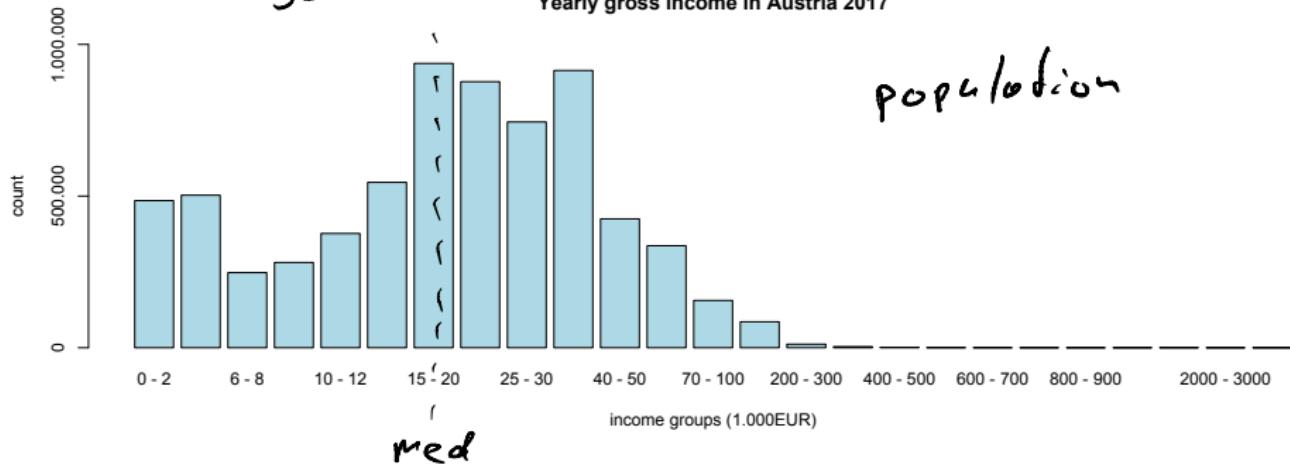
EXAMPLE: MEDIAN INCOME



50% ! 50%

Yearly gross income in Austria 2017

population

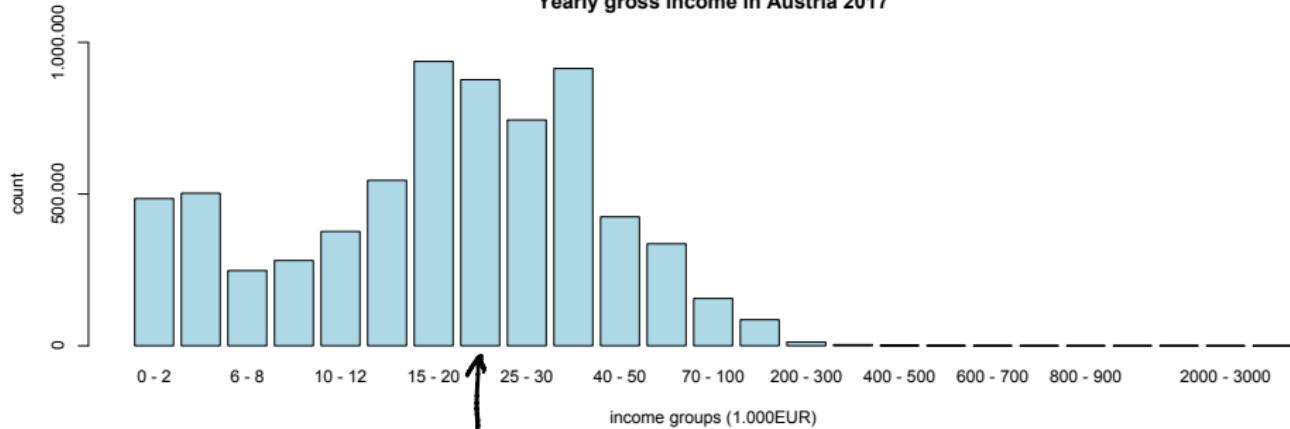


Recall: Median = 50% quantile = smallest number $m \in \mathbb{R}$
such that 50% of data are below or equal to m

EXAMPLE: MEDIAN INCOME



Yearly gross income in Austria 2017



> $N/2$

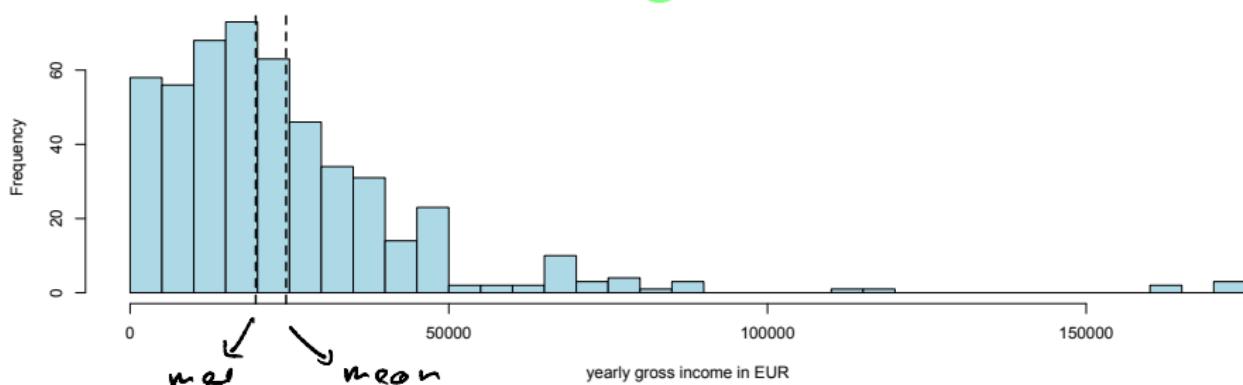
[1] 3466764

> cumsum(rev(dat\$count))

[1] 485050 987971 1235594 1516636 1893307 2438570
[7] 3375906 4252996 4997282 5911340 6336321 6672594
[13] 6828489 6914181 6926012 6929730 6931288 6932039
[19] 6932500 6932774 6932960 6933085 6933416 6933475
[25] 6933529

EXAMPLE: MEDIAN INCOME

Histogram of sample of size 500 of yearly gross income



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
8	10716	<u>19712</u>	24475	31148	174629

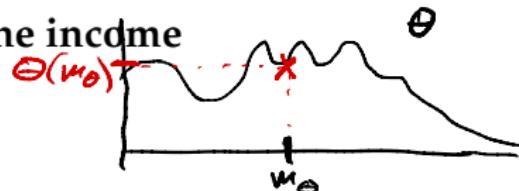
Median should be in the $20k - 25k$ group!?

EXAMPLE: MEDIAN INCOME



universität
wien

A realistic (nonparametric) model for the income distribution:



$$\mathcal{X} = [0, \infty)^n,$$

$$\Theta = \{\text{all pdfs } \theta \in C^1([0, \infty)) : \theta(m_\theta) > 0\}, m_\theta := \text{med}(\theta)$$

$$f_\theta(x) = \prod_{i=1}^n \theta(x_i), \quad x = (x_1, \dots, x_n)' \in \mathcal{X}, \theta \in \Theta.$$

True population is discrete! Should use pmfs?

→ Model is an approximation/idealization!

EXAMPLE: MEDIAN INCOME

A realistic (nonparametric) model for the income distribution:

$$\mathcal{X} = [0, \infty)^n,$$

$$\Theta = \{\text{all pdfs } \theta \in C^1([0, \infty)) : \theta(m_\theta) > 0\}, m_\theta := \text{med}(\theta)$$

$$f_\theta(x) = \prod_{i=1}^n \theta(x_i), \quad x = (x_1, \dots, x_n)' \in \mathcal{X}, \theta \in \Theta.$$

Use sample median $\hat{m}_n := \hat{F}_n^\dagger(1/2)$ to estimate population median $m_\theta := \psi(\theta) := F_\theta^\dagger(1/2)$, $F_\theta(x) := \int_{-\infty}^x \theta(y) dy$.

It is well known that for every $\theta \in \Theta$,

plug-in for θ ?
estimate $\hat{\theta}$
non-parametrically

$$\frac{\hat{m}_n - m_\theta}{se_\theta} \xrightarrow[n \rightarrow \infty]{d.} N(0, 1)$$

where $se_\theta = \frac{1}{2\sqrt{n}\theta(m_\theta)}$.

How do we estimate that?!

BOOTSTRAP ESTIMATION OF SE



Given data $X = (X_1, \dots, X_n)'$, $X_1, \dots, X_n \stackrel{iid}{\sim} p_\theta$, $\theta \in \Theta$, and an estimator $\hat{\theta}_n : \mathcal{X}_0^n \rightarrow \mathbb{R}$ of θ (or $\psi(\theta)$), we want to estimate

$$se_\theta(\hat{\theta}_n) = \sqrt{\text{Var}_\theta[\hat{\theta}_n]}.$$

1. Draw a large number B of random samples of size n (with replacement) from the sample!

$$X_1^* = (X_{1,1}^*, \dots, X_{n,1}^*)'$$

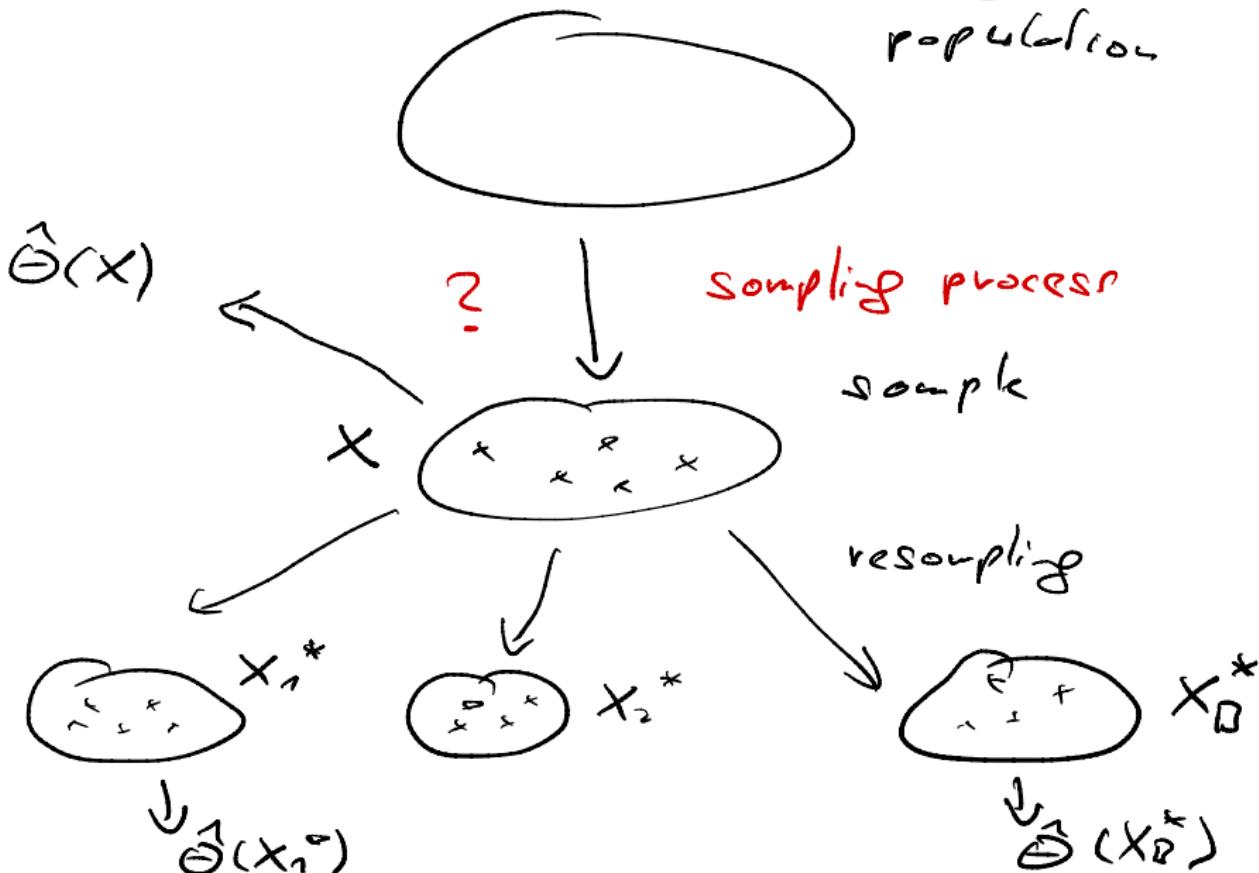
$$X_2^* = (X_{1,2}^*, \dots, X_{n,2}^*)'$$

⋮

$$X_B^* = (X_{1,B}^*, \dots, X_{n,B}^*)'$$

2. Compute $\hat{\theta}_n(X_1^*), \dots, \hat{\theta}_n(X_B^*)$.
3. Approximate $\text{Var}_\theta[\hat{\theta}_n]$ using the LLN (MC idea)

$$\hat{s}e_{boot}^2 = \frac{1}{B} \sum_{j=1}^B \hat{\theta}_n(X_j^*)^2 - \left(\frac{1}{B} \sum_{j=1}^B \hat{\theta}_n(X_j^*) \right)^2.$$





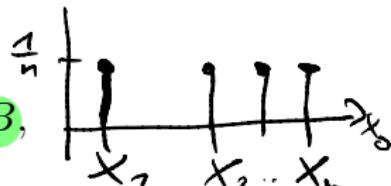
Pretend that the sample is the population and simulate (imitate) the sampling process by 'resampling' from the sample.

Why should this work?

Assume that the original data sample $X = (X_1, \dots, X_n)'$ from sample space $\mathcal{X} = \mathcal{X}_0^n$ is **given and non-random**.

- ▶ Then the **Bootstrap draws**

$$X_{i,j}^*, i = 1, \dots, n, j = 1, \dots, B,$$



are iid from the **empirical distribution** \hat{F}_n of the data, which has pmf

$$\hat{p}_n(x_0) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i\}}(x_0), \quad x_0 \in \{X_1, \dots, X_n\}.$$

- ▶ Thus, also the X_1^*, \dots, X_B^* are iid from the n -fold product

$$\hat{f}_n(x) := \prod_{i=1}^n \hat{p}_n(x_i), \quad x = (x_1, \dots, x_n)' \in \{X_1, \dots, X_n\}^n.$$

- But if $X_1^*, \dots, X_B^* \stackrel{iid}{\sim} \hat{f}_n$, then (by the LLN)

$$\hat{s}e_{boot}^2 = \frac{1}{B} \sum_{j=1}^B \hat{\theta}_n(X_j^*)^2 - \left(\frac{1}{B} \sum_{j=1}^B \hat{\theta}_n(X_j^*) \right)^2 \xrightarrow[B \rightarrow \infty]{i.p.} \text{Var}_{\hat{f}_n}[\hat{\theta}_n].$$

- Because $X_1^* = (X_{1,1}^*, \dots, X_{n,1}^*)' \sim \hat{f}_n$ follows an iid model with marginal pmf \hat{p}_n , we also write $\text{Var}_{\hat{f}_n}[\hat{\theta}_n] = \text{Var}_{\hat{p}_n}[\hat{\theta}_n]$.
- Since the sample $X_1, \dots, X_n \stackrel{iid}{\sim} p_\theta$, we have $\hat{F}_n \approx F_\theta$ (recall HW3), $\hat{p}_n \approx p_\theta$ and one can often show that

$$\text{Var}_{\hat{p}_n}[\hat{\theta}_n] \approx \text{Var}_{p_\theta}[\hat{\theta}_n] = \text{Var}_\theta[\hat{\theta}_n], \quad \text{if } n \text{ is large.}$$

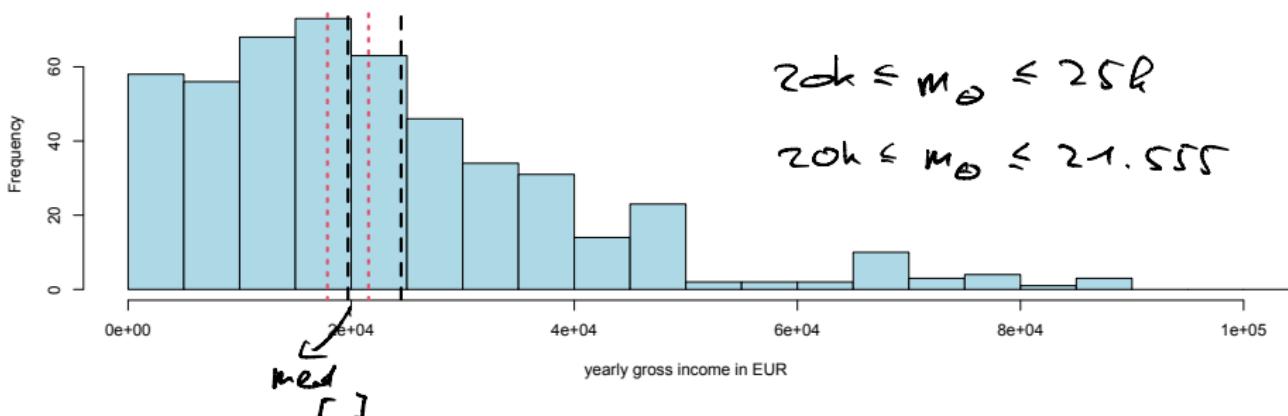
$$\underbrace{\hat{s}e_{boot}^2 \approx \text{Var}_{\hat{p}_n}[\hat{\theta}_n]}_{\text{if } B \text{ is large}}, \quad \underbrace{\text{Var}_{\hat{p}_n}[\hat{\theta}_n] \approx \text{Var}_\theta[\hat{\theta}_n] = se_\theta^2(\hat{\theta}_n)}_{\text{if } n \text{ is large}}$$

EXAMPLE: MEDIAN INCOME

Compute an approximate 95% ($\alpha = 0.05$) bootstrap confidence interval (by bootstrap estimation of the standard error)

$$CI_{\alpha} = [\hat{m}_n - q_{1-\frac{\alpha}{2}}^{(N)} \hat{s}e_{boot}, \hat{m}_n + q_{1-\frac{\alpha}{2}}^{(N)} \hat{s}e_{boot}]$$

Histogram of sample of size 500 of yearly gross income



```
> c(CI_l, CI_u)  
[1] 17870 21555
```

```
> CI_u - CI_l  
[1] 3685
```

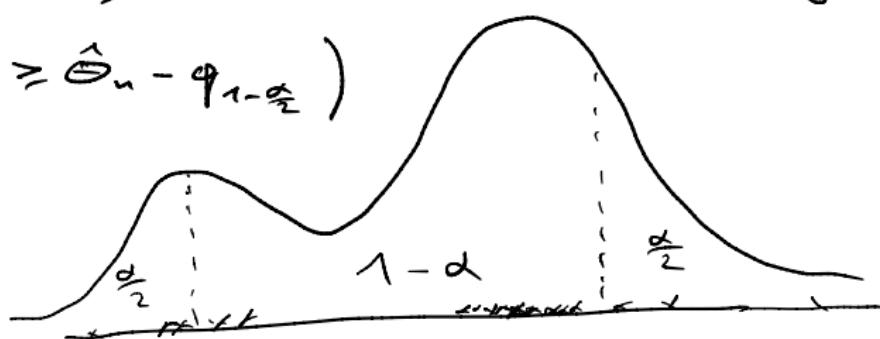
Given data $X = (X_1, \dots, X_n)'$, $X_1, \dots, X_n \stackrel{iid}{\sim} p_\theta$, $\theta \in \Theta$, and an estimator $\hat{\theta}_n : \mathcal{X}_0^n \rightarrow \mathbb{R}$ of θ (or $\psi(\theta)$), we want to estimate the $\alpha/2$ and $1 - \alpha/2$ quantiles of the distribution of $\hat{\theta}_n(X) - \theta$.

Why? Say $\hat{\theta}_n - \theta \sim g_n$ is the unknown sampling distribution.

$$P_\theta(q_{\frac{\alpha}{2}} \leq \hat{\theta}_n - \theta \leq q_{1-\frac{\alpha}{2}}) = 1 - \alpha$$

$$\hat{\theta}_n - \theta \sim g_n$$

$$= P_\theta(\hat{\theta}_n - q_{\frac{\alpha}{2}} \geq \theta \geq \hat{\theta}_n - q_{1-\frac{\alpha}{2}})$$



$$CI_\alpha = [\hat{\theta}_n - q_{1-\frac{\alpha}{2}}, \hat{\theta}_n - q_{\frac{\alpha}{2}}]$$

Given data $X = (X_1, \dots, X_n)'$, $X_1, \dots, X_n \stackrel{iid}{\sim} p_\theta$, $\theta \in \Theta$, and an estimator $\hat{\theta}_n : \mathcal{X}_0^n \rightarrow \mathbb{R}$ of θ (or $\psi(\theta)$), we want to estimate the $\alpha/2$ and $1 - \alpha/2$ quantiles of the distribution of $\hat{\theta}_n(X) - \theta$.

1. Draw a large number B of random samples of size n (with replacement) from the sample!

$$X_1^*, X_2^*, \dots, X_B^*$$

2. Compute empirical (bootstrap) quantiles $\hat{q}_{\alpha/2}^*$ and $\hat{q}_{1-\frac{\alpha}{2}}^*$ of $\hat{\theta}_n(X_1^*) - \hat{\theta}_n(X)$, ..., $\hat{\theta}_n(X_B^*) - \hat{\theta}_n(X)$.
3. $CI_\alpha = [\hat{\theta}_n(X) - \hat{q}_{1-\frac{\alpha}{2}}^*, \hat{\theta}_n(X) - \hat{q}_{\alpha/2}^*]$

EXAMPLE: MEDIAN INCOME



Comparison of 95% bootstrap CIs:

$B = 100$

method	lower	median	upper	length
bootstrap se	17870	19712	21555	3685
pivotal bootstrap	18074	19712	21452	3378

$B = 1000$

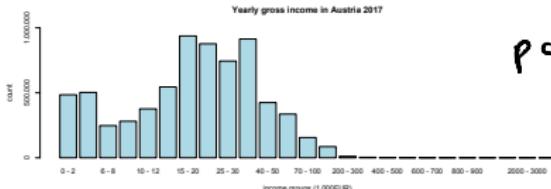
method	lower	median	upper	length
bootstrap se	18100	19712	21324	3224
pivotal bootstrap	18161	19712	21388	3227

EXAMPLE: MEDIAN INCOME



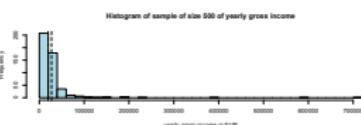
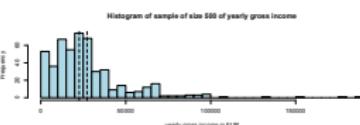
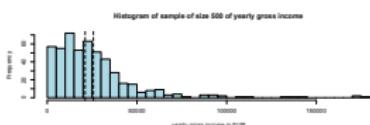
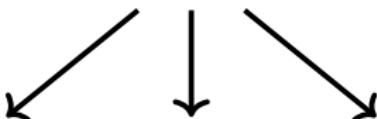
Can we actually trust this CI? Check by simulation!

$$m_\theta = F_\theta^\dagger(1/2)$$

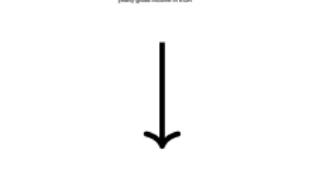


population

samples of size
 $n = 500$



m_θ
Blood shop
sample



EXAMPLE: MEDIAN INCOME



Simulation results:

$$1 - \alpha = 0.95$$

number of (Montecarlo) samples drawn = 1000

$$B = 1000$$

sample size $n = 100$

method	coverage prob.	average length
se bootstrap	0.932	7601
pivotal bootstrap	0.871	7404

$$\geq 1 - \alpha = 0.95$$

sample size $n = 1000$

method	coverage prob.	average length
se bootstrap	0.935	2412
pivotal bootstrap	0.917	2377

Failure of the Bootstrap



- ▶ Regard your data sample as the population.
- ▶ Draw B iid random (bootstrap) samples from the sample (re-sample) like in a MC simulation.
- ▶ Compute your estimator on each of the bootstrap samples.
- ▶ Use this resulting 'bootstrap distribution' of your estimator as an approximation of its true unknown sampling distribution.

- ▶ In the median income example we simulated the actual coverage probability of the bootstrap CIs and found

$$P_\theta(m_\theta \in CI_\alpha) \approx 1 - \alpha.$$

- ▶ We did that for the true data generating parameter $\theta \in \Theta$.
(we cheated!)
- ▶ We concluded that the bootstrap works relatively well for n large.
- ▶ In practice we would try many different choices for θ .



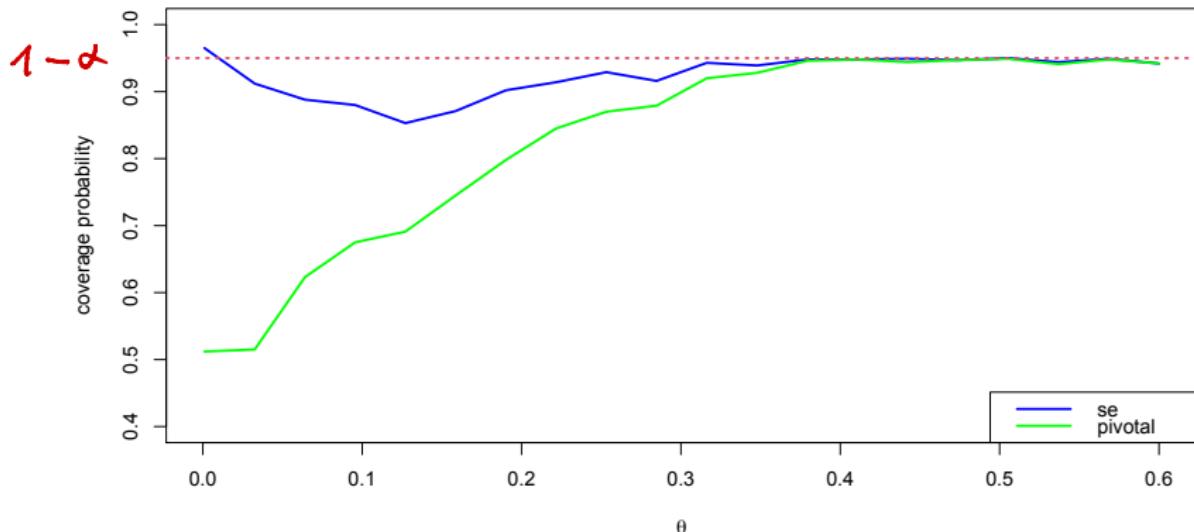
- ▶ Consider $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$, $\theta \in \Theta = [0, \infty)$.
- ▶ Goal: Bootstrap inference on θ .
- ▶ The classical estimator (MLE) is

$$\hat{\theta}_n = \max\{0, \bar{X}_n\}, \quad \bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i.$$

FAILURE OF THE BOOTSTRAP



- ▶ $n = 100, \alpha = 0.05, B = 1000$
- ▶ 1000 MC iterations



FAILURE OF THE BOOTSTRAP

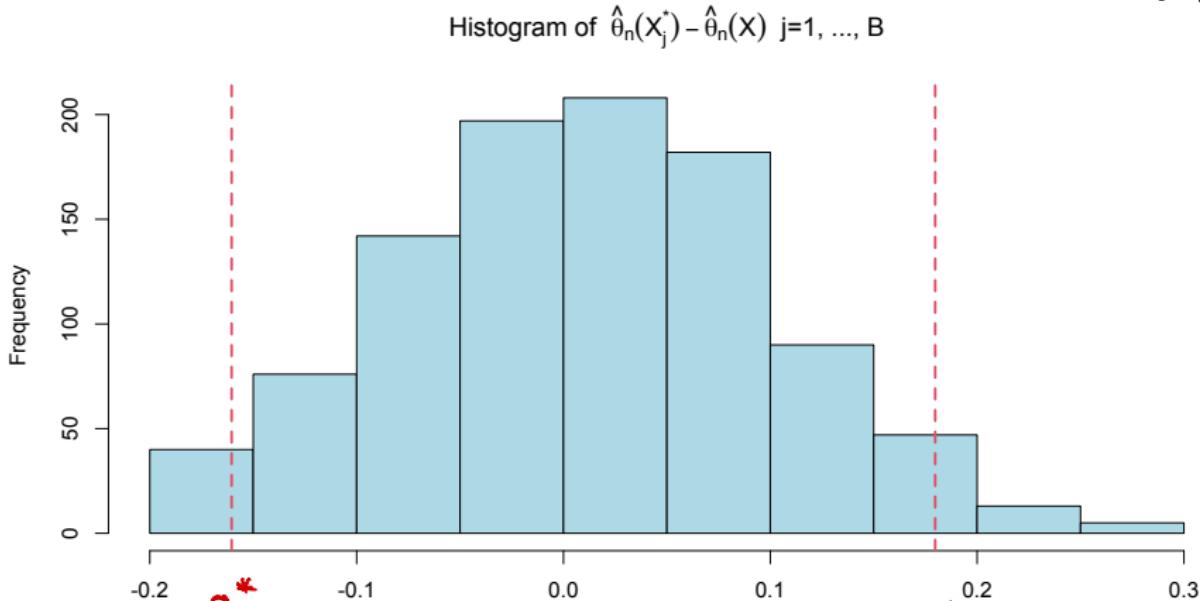


What went wrong?

Look at the 'bootstrap distribution' for a sample

$$X = (X_1, \dots, X_n) \text{ with } \bar{X}_n = 0.185.$$

$$\hat{\theta}_n(x_i^*) = \max\{\theta_n(\bar{X}_j^*)\}$$



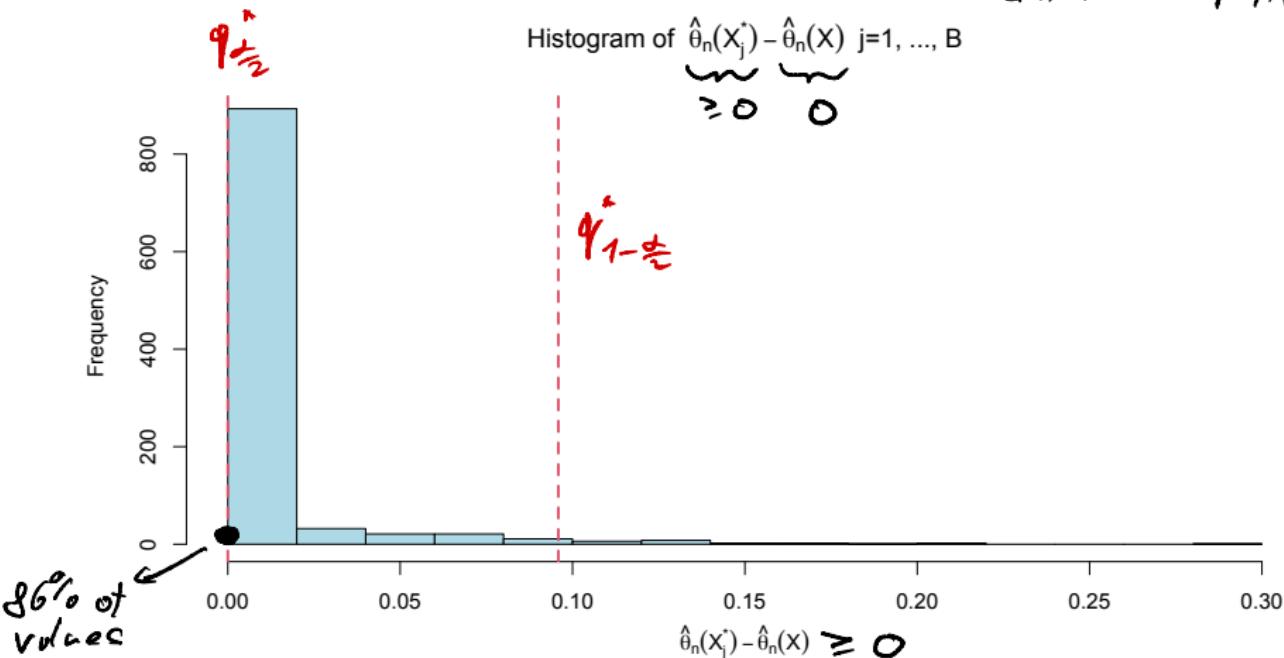
$$\left[\hat{\theta}_n(x) - q_{1-\frac{\alpha}{2}}^*, \hat{\theta}_n(x) - q_{\frac{\alpha}{2}}^* \right]$$

FAILURE OF THE BOOTSTRAP

Look at the 'bootstrap distribution' for a sample with

$$\bar{X}_n = -0.122.$$

$$\hat{\Theta}_-(X) = \max\{0, \hat{X}_-\}$$



Here, actually 86% of bootstrap samples X_j^* produce
 $\hat{\theta}_n(X_j^*) - \hat{\theta}_n(X) = 0.$

- ▶ Look at the ‘bootstrap distribution’ for a sample with $\bar{X}_n = -0.122$.
- ▶ Notice: for every $j \in \{1, \dots, B\}$, we have

$$\hat{\theta}_n(X_j^*) - \hat{\theta}_n(X) = \max\{0, \bar{X}_j^*\} - \max\{0, \bar{X}_n\} \geq 0,$$
$$\Rightarrow 0 \leq \hat{q}_{\alpha/2}^* \leq \hat{q}_{1-\alpha/2}^*.$$

- ▶ Therefore,

$$CI_\alpha = [\underbrace{\hat{\theta}_n(X) - \hat{q}_{1-\alpha/2}^*}_{=0}, \underbrace{\hat{\theta}_n(X) - \hat{q}_{\alpha/2}^*}_{=0}] \subseteq (-\infty, 0] \quad \forall \theta > 0.$$

FAILURE OF THE BOOTSTRAP



Recall: $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$

$\alpha = 0, 05$

Thus, intuitively, for very small $\theta > 0$,

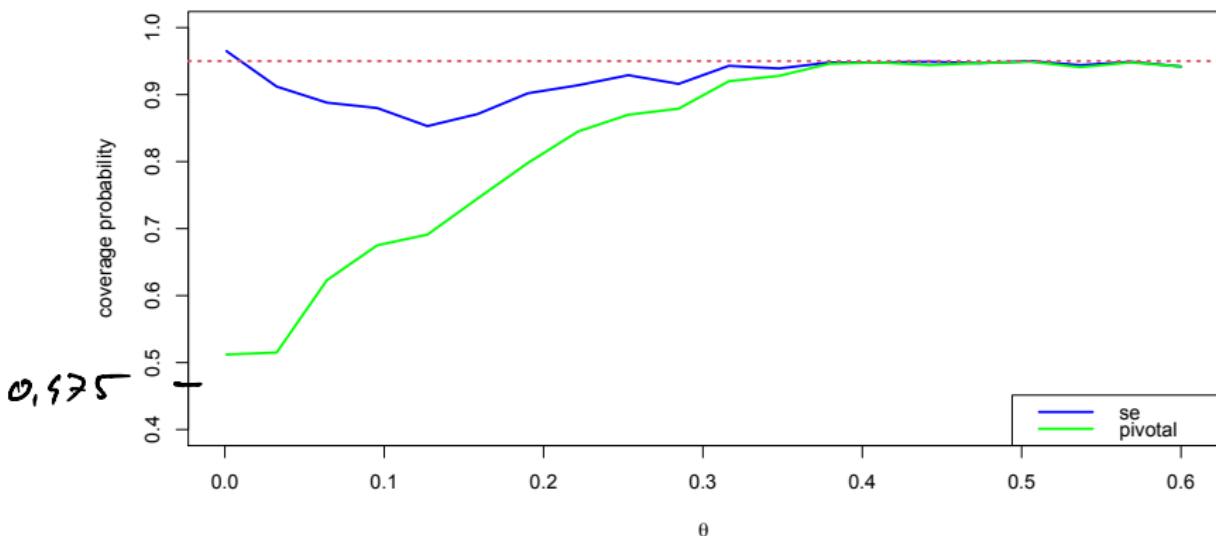
$$P_\theta(\theta \in CI_{0.05}) = P_\theta(\theta \in CI_{0.05} | \bar{X}_n < 0)P_\theta(\bar{X}_n < 0) \\ = 0 \quad \approx \frac{1}{2}$$

$$\bar{X}_n \sim N(\theta, \frac{1}{n}) \quad + P_\theta(\theta \in CI_{0.05} | \bar{X}_n \geq 0)P_\theta(\bar{X}_n \geq 0) \\ \approx 0, 95 \quad \leq 1 \quad \approx \frac{1}{2} \\ \approx 0, 475$$

FAILURE OF THE BOOTSTRAP



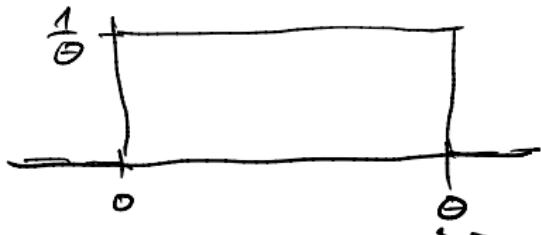
- ▶ $n = 100, \alpha = 0.05, B = 1000$
- ▶ 1000 MC iterations



Is the se-bootstrap always superior to the pivotal method?

- ▶ Consider $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}[0, \theta]$, $\theta \in \Theta = (0, \infty)$.
- ▶ Goal: Bootstrap inference on θ .
- ▶ The classical estimator (MLE) is

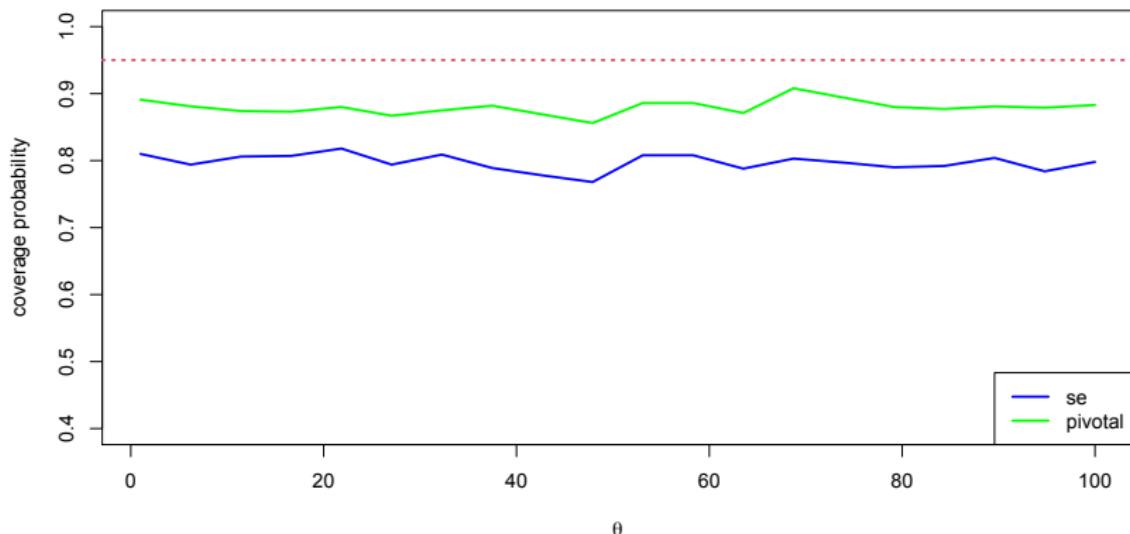
$$\hat{\theta}_n = \max\{X_1, \dots, X_n\}.$$



FAILURE OF THE BOOTSTRAP



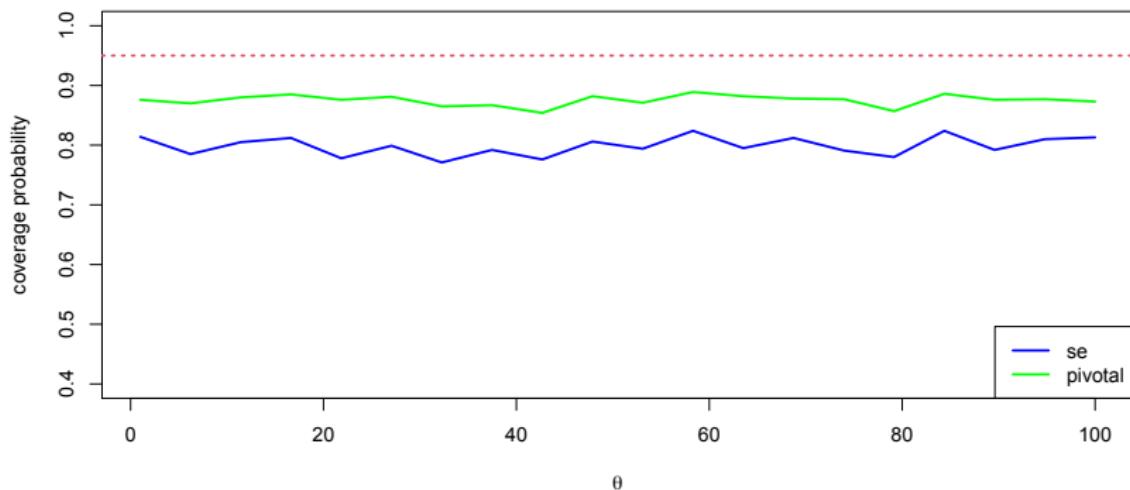
- ▶ $n = 100, \alpha = 0.05, B = 1000$
- ▶ 1000 MC iterations



FAILURE OF THE BOOTSTRAP



- ▶ $n = 1000, \alpha = 0.05, B = 1000$
- ▶ 1000 MC iterations

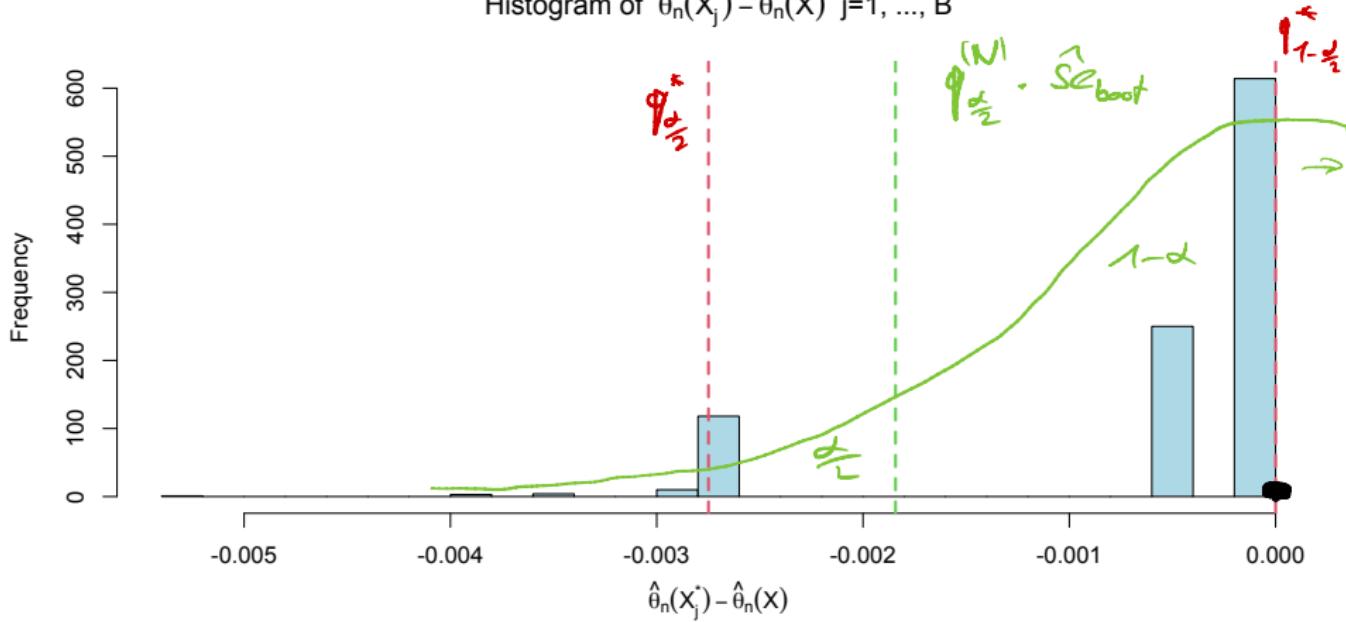


FAILURE OF THE BOOTSTRAP



Look at the bootstrap distribution of one given sample of size $n = 1000$ with $\theta = 1$.

Histogram of $\hat{\theta}_n(X_j^*) - \hat{\theta}_n(X)$ $j=1, \dots, B$



Here, 62% of all bootstrap samples X_j^* produce $\hat{\theta}_n(X_j^*) - \hat{\theta}_n(X) = 0$.

Here, 62% of all bootstrap samples X_j^* produce

$$\hat{\theta}_n(X_j^*) - \hat{\theta}_n(X) = 0.$$

Why?
assume no ties

$$\hat{\theta}_n(X_1^*) - \hat{\theta}_n(X) = \max\{X_{1,1}^*, \dots, X_{n,1}^*\} - \max\{X_1, \dots, X_n\}$$

- ▶ This is equal to 0 if, and only if, in our bootstrap sample $X_{1,1}^*, \dots, X_{n,1}^*$ we happen to draw the largest sample point from $\{X_1, \dots, X_n\}$.
- ▶ What is the probability of that happening?

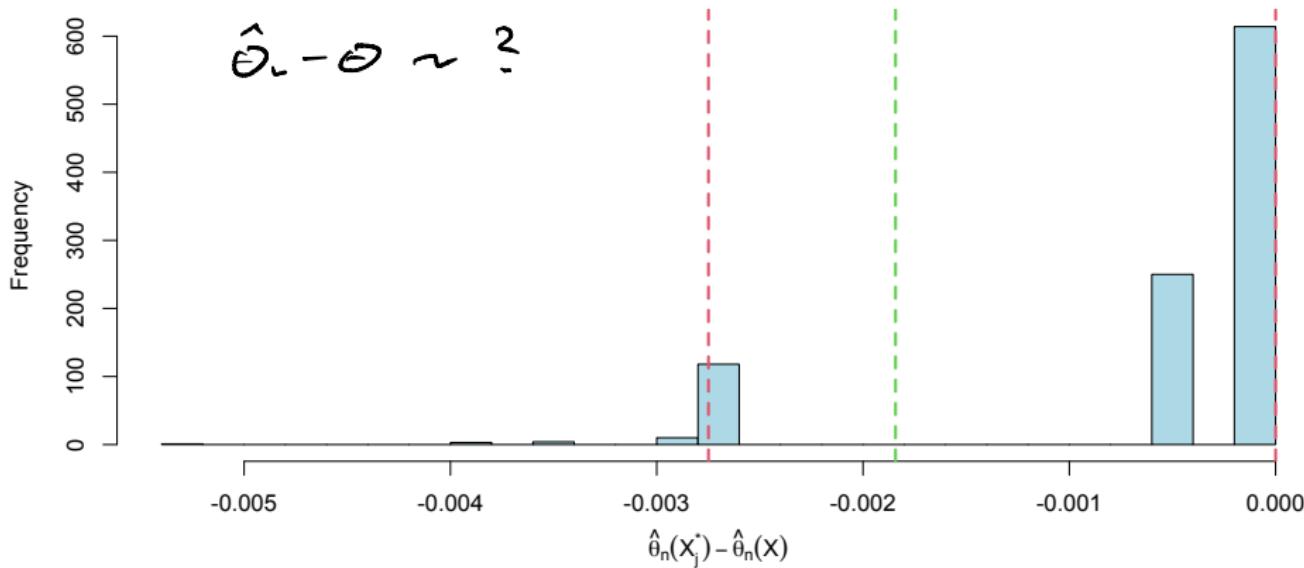
$$1 - \left(\frac{n-1}{n}\right)^n = 1 - \left(1 - \frac{1}{n}\right)^n \xrightarrow{n \rightarrow \infty} 1 - e^{-1} \approx 0,63$$

FAILURE OF THE BOOTSTRAP



Look at the bootstrap distribution of one given sample of size $n = 1000$ with $\theta = 1$.

Histogram of $\hat{\theta}_n(X_j^*) - \hat{\theta}_n(X)$ $j=1, \dots, B$



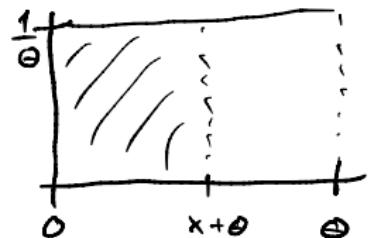
Here, 62% of all bootstrap samples X_j^* produce $\hat{\theta}_n(X_j^*) - \hat{\theta}_n(X) = 0$.

What is the actual sampling distribution of $\hat{\theta}_n - \theta$?

Recall:

$$X_1, \dots, X_n \stackrel{iid}{\sim} \text{Unif}[0, \theta], \theta \in \Theta = (0, \infty), \hat{\theta}_n = \max\{X_1, \dots, X_n\}.$$

$$\begin{aligned} P_\theta(\hat{\theta}_n - \theta \leq x) &= P_\theta\left(\max_{1 \leq i \leq n} \{X_i\} \leq x + \theta\right) \\ &= P_\theta(X_1 \leq x + \theta \text{ and } \dots \text{ and } X_n \leq x + \theta) \\ &\stackrel{\text{indep.}}{=} \prod_{i=1}^n P_\theta(X_i \leq x + \theta) \end{aligned}$$



$$= \begin{cases} 0 & \text{if } x + \theta < 0, \text{ if } x < -\theta \\ \frac{x+\theta}{\theta} & \text{if } 0 \leq x + \theta \leq \theta, -\theta \leq x \leq 0 \\ 1 & \text{if } x + \theta > \theta, x > 0 \end{cases}$$

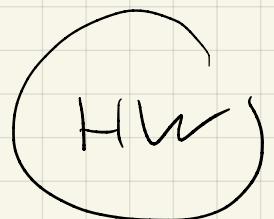
$$P_0(\bar{\Theta}_n - \Theta \leq x) = \begin{cases} 0 & x < -\Theta \\ \left(1 + \frac{x}{\Theta}\right)^n & -\Theta \leq x \leq 0 \\ 1 & x > 0 \end{cases}$$

$$g_n(x) = \frac{d}{dx} P_0(\bar{\Theta}_n - \Theta \leq x)$$

$$= \begin{cases} 0 & x < -\Theta \\ n \left(1 + \frac{x}{\Theta}\right)^{n-1} \frac{1}{\Theta}, & -\Theta \leq x \leq 0 \\ 0, & x > 0 \end{cases}$$

pdf of sampling distribution

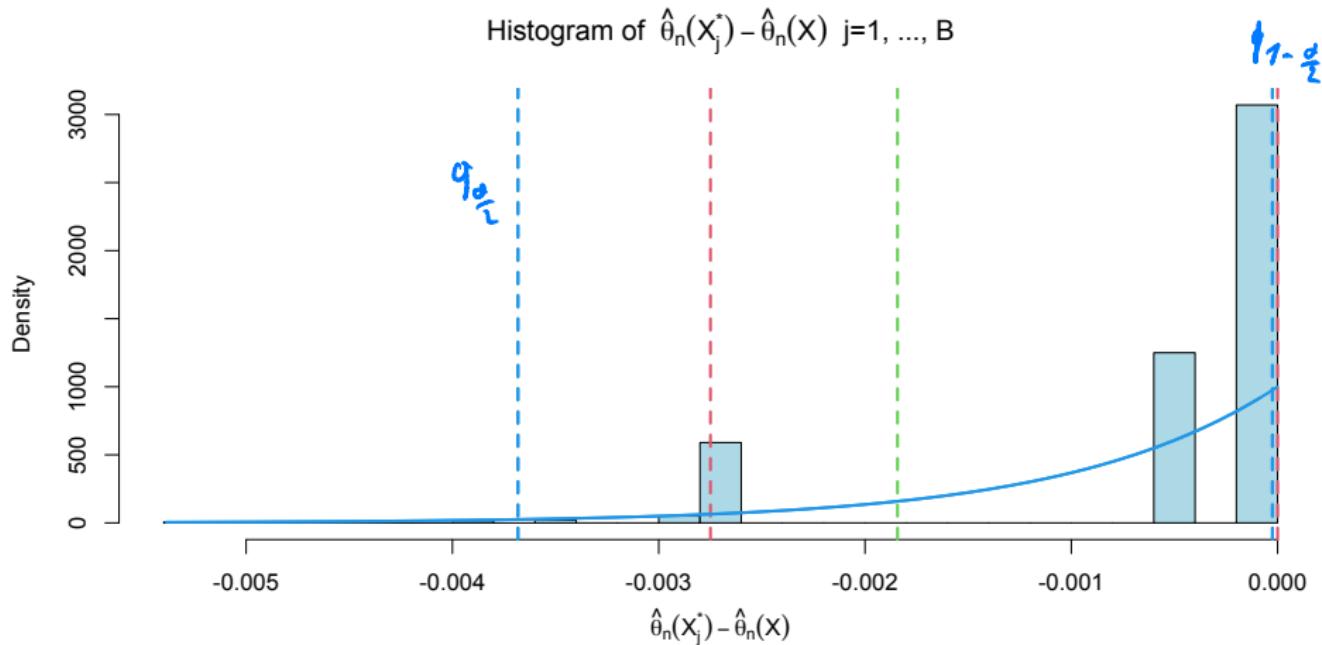
$$f_x(g_n) = \left[x^{\frac{1}{n}} - 1 \right] \Theta$$



FAILURE OF THE BOOTSTRAP



Bootstrap vs. true sampling distribution of $\hat{\theta}_n - \theta$. $\theta = 1, n = 1000.$



Here, 62% of all bootstrap samples X_j^* produce $\hat{\theta}_n(X_j^*) - \hat{\theta}_n(X) = 0.$



**Always take a look at the
bootstrap distribution!!!**



Summing up:

- ▶ Computationally expensive simulation method.
- ▶ Very flexible, generic, no need for problem specific formulas.
- ▶ Produces only approximate inference; requires large samples.
- ▶ Can go wrong!
- ▶ In general: Works, if the normal approximation works.