# STA 104 Applied Nonparametric Statistics

## Chapter 6: Bootstrap

Xiner Zhou

Department of Statistics, University of California, Davis

# Table of contents

# Bootstrap Confidence Intervals

$$P\left( L(data) < \theta < U(data) \right) = \underbrace{1 - \alpha}_{95\%} \quad . \quad \alpha = 4.5$$

fixed (pointing to $\theta$)

random (pointing to $L(data)$ and $U(data)$)

$\underbrace{\phantom{P(L(data) < \theta < U(data))}}$ defines a $100(1-\alpha)\%$ C.I. for $\theta$

a region of values, depending on data,
s.t. true $\theta$ lies within this interval
is $100(1-\alpha)\%$.

unknown population

$\theta$: typical value

$\theta_0$

$\theta$

$\theta$

$\theta$

$L(data)$     $U(data)$

$\theta$

# Example 1: Confidence Interval for $\mu$

Suppose that $x_i \overset{iid}{\sim} N\left(\mu, \sigma^2\right)$ for $i = 1, \ldots, n$ and we want to form a confidence interval for $\mu$. As an estimate of $\mu$, we will use the sample mean $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$. Assuming that $x_i \overset{iid}{\sim} N\left(\mu, \sigma^2\right)$, we know that $\bar{x} \sim N\left(\mu, \sigma^2/n\right)$, which implies that $\sqrt{n}(\bar{x} - \mu)/\sigma \sim N(0, 1)$. As a result, we have that

*simply distr*

$$P\left(z_{\alpha/2} < \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} < z_{1-\alpha/2}\right) = 1 - \alpha$$

where $z_\alpha = \Phi^{-1}(\alpha)$ with $\Phi^{-1}(\cdot)$ denoting the quantile function for the standard normal distribution. Rearranging the terms inside the above probability statement gives

$$1 - \alpha = P\left(z_{\alpha/2}\sigma/\sqrt{n} < \bar{x} - \mu < z_{1-\alpha/2}\sigma/\sqrt{n}\right)$$
$$= P\left(z_{\alpha/2}\sigma/\sqrt{n} - \bar{x} < -\mu < z_{1-\alpha/2}\sigma/\sqrt{n} - \bar{x}\right)$$
$$= P\left(\bar{x} + z_{\alpha/2}\sigma/\sqrt{n} > \mu > \bar{x} - z_{1-\alpha/2}\sigma/\sqrt{n}\right)$$

which implies that a $100(1 - \alpha)\%$ confidence interval for $\mu$ defines $a(\bar{x}) = \bar{x} - z_{1-\alpha/2}\sigma/\sqrt{n}$ and $b(\bar{x}) = \bar{x} - z_{\alpha/2}\sigma/\sqrt{n}$. Note that since $-z_{\alpha/2} = z_{1-\alpha/2}$ we can write the two endpoints of the confidence interval as

$$\bar{x} \pm z_{1-\alpha/2}\text{SE}(\bar{x})$$

where $\text{SE}(\bar{x}) = \sigma/\sqrt{n}$ is the standard error of the sample mean. In practice, it is typical to form a 90% confidence interval (i.e., $\alpha = 0.1$), which corresponds to $z_{0.95} \approx 1.65$, a 95% confidence interval (i.e., $\alpha = 0.05$), which corresponds to $z_{0.975} \approx 1.96$, or a 99% confidence interval (i.e., $\alpha = 0.01$), which corresponds to $z_{0.995} = 2.58$.

Forming a confidence interval for $\mu$ with $x_i \overset{\text{iid}}{\sim} N(\mu, \sigma^2)$, here is a simple demonstration of forming a 95% confidence interval using $R = 10000$ replications with $n = 25$ observations.

```
 R <- 10000
n <- 25
set.seed(1)
xbar <- replicate(R, mean(rnorm(n)))
ci.lo <- xbar - qnorm(.975) / sqrt(n)    # 95% CI lower bound
ci.up <- xbar - qnorm(.025) / sqrt(n)    # 95% CI upper bound
ci.in <- (ci.lo <= 0) & (0 <= ci.up)   ←  check  0 ∈ CI.
mean(ci.in)

## [1] 0.9501
```

# Example 2: Confidence Interval for $\sigma^2$

Suppose that $x_i \overset{\text{iid}}{\sim} N\left(\mu, \sigma^2\right)$ for $i = 1, \ldots, n$ and we want to form a confidence interval for $\sigma^2$. As an estimate of $\sigma^2$, we will use the sample variance $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$. Assuming that $x_i \overset{\text{iid}}{\sim} N\left(\mu, \sigma^2\right)$, we know that $(n-1)s^2/\sigma^2 \sim \chi^2_{n-1}$.

As a result, we have that

*need to know*

$$P\left(q_{n-1;\alpha/2} < (n-1)\frac{s^2}{\sigma^2} < q_{n-1;1-\alpha/2}\right) = 1 - \alpha$$
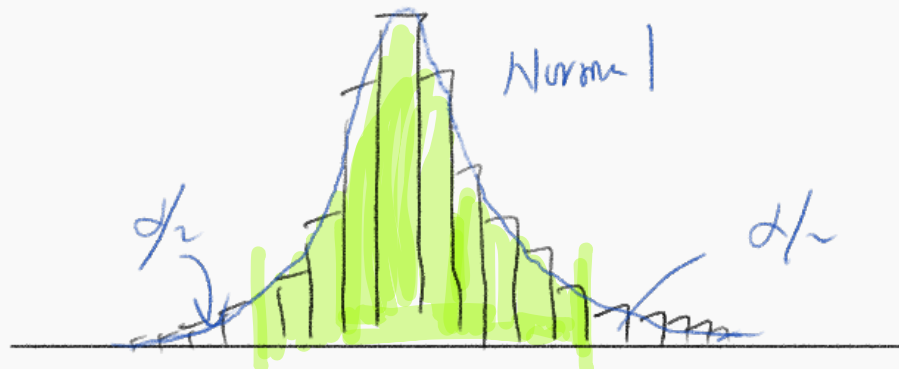
where $q_{n-1;\alpha} = Q_{n-1}(\alpha)$ with $Q_{n-1}(\cdot)$ denoting the quantile function for the $\chi^2_{n-1}$ distribution. Rearranging the terms inside the above probability statement gives

$$1 - \alpha = P\left(\frac{q_{n-1;\alpha/2}}{n-1} < \frac{s^2}{\sigma^2} < \frac{q_{n-1;1-\alpha/2}}{n-1}\right)$$

$$= P\left(\frac{q_{n-1;\alpha/2}}{s^2(n-1)} < \frac{1}{\sigma^2} < \frac{q_{n-1;1-\alpha/2}}{s^2(n-1)}\right)$$

$$= P\left(\frac{s^2(n-1)}{q_{n-1;\alpha/2}} > \sigma^2 > \frac{s^2(n-1)}{q_{n-1;1-\alpha/2}}\right)$$

*here!*

which implies that a $100(1-\alpha)\%$ confidence interval for $\sigma^2$ defines $a\left(s^2\right) = (n-1)s^2/q_{n-1;1-\alpha/2}$ and $b\left(s^2\right) = (n-1)s^2/q_{n-1;\alpha/2}$.
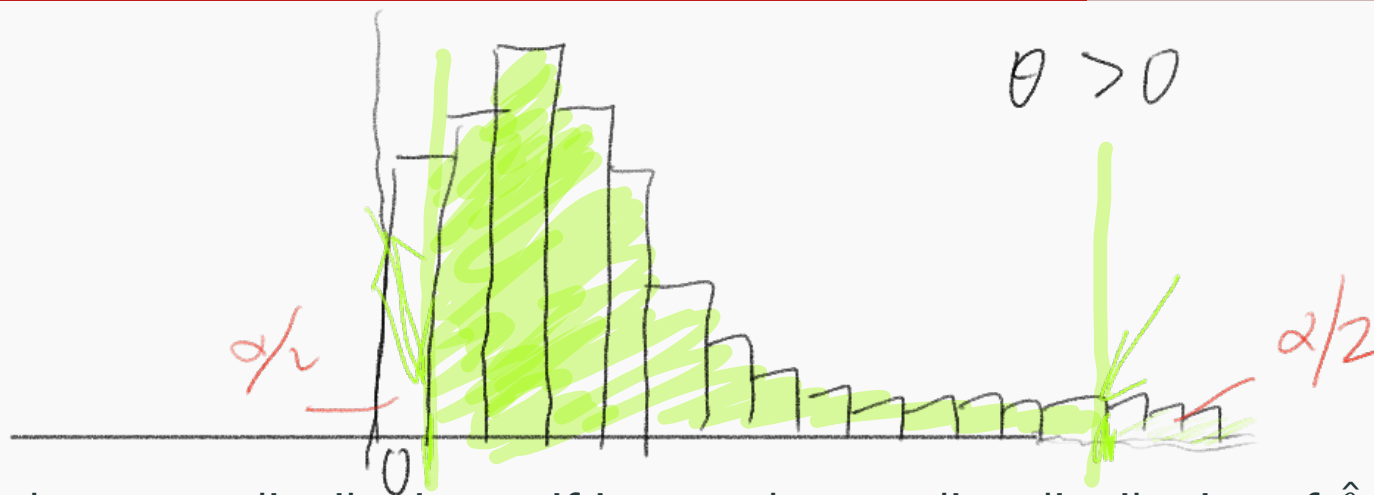
# Bootstrap CI: Normal Approximation



The normal approximation confidence interval uses the classic confidence interval formula (for the mean), but replaces the standard error with the bootstrap estimate of the standard error.

Specifically, the normal approximation interval has the form

$$\hat{\theta} \pm Z_{1-\alpha/2} \widehat{SE}(\hat{\theta})$$

where $Z_{1-\alpha/2}$ is the quantile of the standard normal distribution that cuts-off $\alpha/2$ in the upper tail (e.g., $Z_{1-\alpha/2} = 1.96$ for a 95% interval), and $\widehat{SE}(\hat{\theta})$ is the bootstrap estimate of the standard error of $\hat{\theta}$.

# Bootstrap CI: Percentile Method



Simply uses the bootstrap distribution as if it were the sampling distribution of $\hat{\theta}$.

The percentile method defines the $100(1-\alpha)\%$ confidence interval for $\theta$ as

$$\left[ Q^*_{\alpha/2}, Q^*_{1-\alpha/2} \right]$$

where $Q^*_{\alpha/2}$ and $Q^*_{1-\alpha/2}$ denote the quantiles of the bootstrap distribution of $\hat{\theta}$.

# Example: To estimate the correlation between Petal Length and Petal Width

```r
library(boot)

# Custom function to find correlation
# between the Petal Length and Width
corr.fun <- function(data, idx)
{
  df <- data[idx, ]

  # Find the spearman correlation between
  # the 3rd and 4th columns of dataset
  c(cor(df[, 3], df[, 4], method = 'spearman'))
}

bootstrap <- boot(iris, corr.fun, R = 1000)
bootstrap

##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
##
## Call:
## boot(data = iris, statistic = corr.fun, R = 1000)
##
##
## Bootstrap Statistics :
##      original       bias    std. error
## t1* 0.9376668 -0.00274143  0.00980653
```

*(handwritten annotations:)* index of re-samples — data frame — resample — orig data

*spearman correlation*

```r
# bootstrap distribution
hist(bootstrap$t, xlab = "Statistic", main = "Bootstrap Distribution")
box()
abline(v = bootstrap$t0, lty = 2, col = "red")
legend("topleft", "t0", lty = 2, col = "red", bty = "n")
```

*observed*

**Bootstrap Distribution**

```r
# Function to find the bootstrap Confidence Intervals
boot.ci(boot.out = bootstrap,
        type = c("norm",
                 "perc", "bca"))
```

```
## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 1000 bootstrap replicates
##
## CALL :
## boot.ci(boot.out = bootstrap, type = c("norm", "perc", "bca"))
##
## Intervals :
## Level      Normal              Percentile           BCa
## 95%   ( 0.9212,  0.9596 )   ( 0.9126,  0.9522 )   ( 0.9173,  0.9539 )
## Calculations and Intervals on Original Scale
```

# Example: To estimate Median

We will generate $n = 100$ observations from a standard normal distribution, and use the median as the parameter/statistic of interest.

```
library(boot)          ← package

# generate 100 standard normal observations
set.seed(1)
n <- 100                      ] generate data
x <- rnorm(n)
sim.data=data.frame(x=x)

median.fun <- function(data, idx)
{                    △    △
  df <- data[idx, ]  ←                      ] calculate median θ̂ for each Bootstrap samples.
  quantile(df,prob=0.5) ←─ median
}

bootstrap <- boot(sim.data, median.fun, R = 1000) ←
bootstrap    △      ↑              ↑
                  original data   5000

##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
##
## Call:
## boot(data = sim.data, statistic = median.fun, R = 1000)
##
##
## Bootstrap Statistics :          ŜE_boot = SE(θ̂)
##        original      bias      std. error
## t1*  0.1139092  0.02243135   0.1413482
```

$$\theta_0 = 0 \cdot \hat{\theta} = 0.11$$

```
# bootstrap distribution
hist(bootstrap$t, xlab = "Statistic", main = "Bootstrap Distribution")
box()
abline(v = bootstrap$t0, lty = 2, col = "red")
legend("topleft", "t0", lty = 2, col = "red", bty = "n")
```



Bootstrap Distribution

*unzip object from boot( )*

```r
# Function to find the bootstrap Confidence Intervals
boot.ci(boot.out = bootstrap,
        type = c("norm",
                 "perc", "bca"))
```

*norm < percentile < bca*

```
## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 1000 bootstrap replicates
##
## CALL :
## boot.ci(boot.out = bootstrap, type = c("norm", "perc", "bca"))
##
## Intervals :
## Level      Normal           Percentile          BCa
## 95%    (-0.1856,  0.3685 )  (-0.0593,  0.3788 )  (-0.0811,  0.3692 )
## Calculations and Intervals on Original Scale
```