

The Bootstrap (I)

Suppose that X_1, \dots, X_n is a random sample from some distribution.

We can estimate the population mean μ by \bar{X} , and the standard error (estimated standard deviation)

is s/\sqrt{n} , where $[s$ is the sample standard deviation.]

Similarly, if we want to estimate a proportion p using n independent trials, the estimate is \hat{p} , with

standard error $SE(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$.

The estimate and the SE can be combined to give

asymptotic confidence intervals as estimate $\pm z_{\alpha/2} SE$.

⇒ But what do we do for other parameters?

比如 median 和 quantiles 而不是 常
见的 mean, proportion

The Bootstrap (II)

定义:

★ ⇒ The bootstrap is a way to find an SE for any estimate or a confidence interval for any parameter.

It doesn't work well in every situation, but tends to work well in many situations, especially those where asymptotic normality holds.

⇒ Key idea: Learn about the variability in an estimate over repeated sampling by using sampling from the sample to approximate sampling from the population.

Why "bootstrap"?: It conveys the idea that we're getting extra information that shouldn't be there in some sense by using the sample in more than one way.

Bootstrap Samples

⇒ Given X_1, \dots, X_n , a bootstrap sample is a sample of size n (drawn with replacement) from the sample.

We get a bootstrap sample by sampling from the distribution w/ distribution function \hat{F} (the EMF).

⇒ Note: This is the nonparametric bootstrap. It is also possible to do a parametric bootstrap by assuming a particular parametric family like the normal family.

★ Ex: If the data are $X_1=1, X_2=4$, list all possible bootstrap samples and their probabilities.

Solution: $2^2 = 4$

<u>Sample</u>	<u>Prob.</u>	<u>Sample</u>	<u>Prob</u>
1, 1	1/4	4, 1	1/4
1, 4	1/4	4, 4	1/4

← ↑
Ordered samples

Small Bootstrap Example

Ex: Find the bootstrap distribution of the sample median if the data are $X_1 = 3, X_2 = 5, X_3 = 7$. ($n = 3$).

Solution: There are $3^3 = 27$ possible samples.

	median		median		median
333	3	533	3	733	3
335	3	535	5	735	5
337	3	537	5	737	7
353	3	553	5	753	5
355	5	555	5	755	5
357	5	557	5	757	7
373	3	573	5	773	7
375	5	575	5	775	7
377	7	577	7	777	7

sample median	prob.
3	$7/27$
5	$13/27$
7	$7/27$

The Bootstrap Via Simulation (I)

When the sample is larger, listing is not feasible.

Instead, we do a simulation study.

Suppose we are estimating a parameter θ using $\hat{\theta}$.

our estimator

Procedure: ① Draw B different bootstrap samples for

$B = 500, 1000, \text{ or more.}$

② Compute the bootstrap estimators $\hat{\theta}_{b,1}, \dots, \hat{\theta}_{b,B}$.

bootstrap $\hat{\theta}$ is sample mean. you
draw B sample mean for
all these bootstrap
example

③ Obtain the estimated MSE (mean squared error)

$$\text{as } \hat{MSE} = \frac{1}{B} \sum_{i=1}^B (\hat{\theta}_{b,i} - \hat{\theta})^2 \text{ and the}$$

(3.4.6)
note: $\hat{\theta}$ plays the role of true parameter value
for ex: the mean for the sample you got

SE as $\sqrt{\hat{MSE}}$.

Estimator $E[(\hat{\theta} - \theta)^2]$.

The Bootstrap Via Simulation (II)

We can also estimate other quantities.

= 0 bias is unbiased estimator

Bias: The bias is $B(\hat{\theta}) = E[\hat{\theta} - \theta] = E[\hat{\theta}] - \theta$.

Estimate this as
$$\hat{B}(\hat{\theta}) = \underbrace{\frac{1}{B} \sum_{b=1}^B \hat{\theta}_{b,i}}_{\hat{E}} - \hat{\theta}$$

plays the role of θ .

为实际参数

Variance: Estimate this as

$$\hat{V}(\hat{\theta}) = \frac{1}{B} \sum_{b=1}^B (\hat{\theta}_{b,i} - \hat{E})^2$$

estimated expected value.

Note: Some references use $\sqrt{\hat{V}(\hat{\theta})}$ as the bootstrap SE. Using $\sqrt{\hat{MSE}(\hat{\theta})}$ allows us to account for bias, too. ($MSE = \text{Variance} + \text{Bias}^2$).

Bootstrap Example

EX: Consider the data $3, 4, 4, 5, 7, 9, 11, 13, 14, 19$.

$n=10$

目的:

Using the bootstrap with $B=1000$ bootstrap samples, estimate the MSE, the bias, + the variance associated with

① Using the sample median to estimate the population median and

② Using the sample standard deviation to estimate the population standard deviation.

In each case, also give the point estimate and SE.

Bootstrap Percentile Confidence Intervals

⇒ One of the simplest methods for obtaining bootstrap confidence intervals is called the percentile method.

We find the B bootstrap point estimates $\hat{\theta}_{b,1}, \dots,$

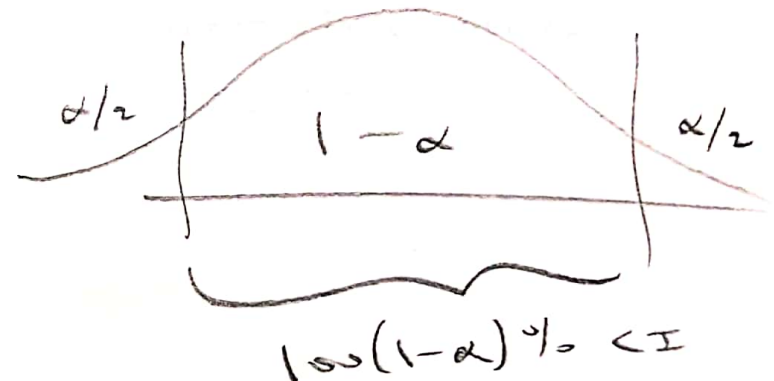
$\hat{\theta}_{b,B}$. We then put these point estimates in

order, and our $100(1-\alpha)\%$ confidence interval for

θ goes from the $(\alpha/2)$ quantile of the bootstrap

$\hat{\theta}$ values to the $1-\alpha/2$ quantile.

Note: These CIs do not necessarily perform well, particularly if the distribution for $\hat{\theta}$ is skewed.



Percentile Confidence Interval Example

Using the 10 data values from earlier, find
95% bootstrap percentile confidence interval
for the population mean, median, and
standard deviation.

One helpful R function: `quantile`.

Using `quantile(X, c(.025, .975))` gives a 95%

↑
bootstrap $\hat{\theta}$
values

← quantiles of interest
bootstrap percentile
confidence interval.

How well do the intervals work?

→ { coverage probability
short CI

Assuming normal data and using $B=1000$,

$n=5, 10, 20$, and confidence level 95%,

estimate the true coverage probability for 95%

bootstrap percentile confidence intervals for the population mean.

Then repeat the same study using exponential data to assess the impact of skewness.

$n=5$ 79%

$n=10$ 86%

$n=20$ 91%

$n=40$ 91.4%

min sample size: $n=5 \Rightarrow 0.761$

for normal data $n=10 \Rightarrow 0.9$

$n=20 \Rightarrow 0.94$

$n=40 \Rightarrow 0.934$

} close 95%

not so great for small sample size

Better CIs by Pivoting

A pivotal quantity is a random variable with a distribution that doesn't depend on any unknown parameters.

Ex: If $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$, then both

① $t = \frac{\overset{\text{Sample mean}}{\bar{X}} - \mu}{s/\sqrt{n}}$ and ② $\chi^2 = \frac{(n-1)s^2}{\sigma^2}$

are pivotal quantities. \because No matter what μ is and σ^2 is

Their distributions?

① t w/ $n-1$ df

② $\chi^2(n-1)$

free of μ and σ^2

Pivotal Quantities (Lead to CI) (I)

Ex: Show how to derive a CI for a normal mean μ by using the fact that $t \sim t_{n-1}$.

Soln: $P(-t_{\alpha/2} \leq \frac{\bar{X} - \mu}{s/\sqrt{n}} \leq t_{\alpha/2}) = 1 - \alpha$

$$\Rightarrow P(\bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}}) = 1 - \alpha$$

$$\Rightarrow \boxed{\bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}} \text{ is a } 100(1-\alpha)\% \text{ CI for } \mu.}$$

Pivotal Quantities Lead to CIs (II)

Ex: Show how to find a CI for a normal variance σ^2 by using the fact that $X^2 \sim \chi^2_{n-1}$.

Soln: $P\left(\chi^2_{n-1, 1-\alpha/2} \leq \frac{(n-1)s^2}{\sigma^2} \leq \chi^2_{n-1, \alpha/2}\right) = 1-\alpha$ In upper tail.

$$\Rightarrow P\left(\frac{1}{\chi^2_{n-1, \alpha/2}} \leq \frac{\sigma^2}{(n-1)s^2} \leq \frac{1}{\chi^2_{n-1, 1-\alpha/2}}\right) = 1-\alpha$$

$$\Rightarrow P\left(\frac{(n-1)s^2}{\chi^2_{n-1, \alpha/2}} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi^2_{n-1, 1-\alpha/2}}\right) = 1-\alpha.$$

$$\Rightarrow \left(\frac{(n-1)s^2}{\chi^2_{n-1, \alpha/2}}, \frac{(n-1)s^2}{\chi^2_{n-1, 1-\alpha/2}} \right) \text{ is a } 100(1-\alpha)\% \text{ CI for } \sigma^2.$$

Bootstrap Pivotal CIs

The idea is to use the bootstrap to estimate the appropriate critical values for non-normal data.

Suppose that the data are non-normal, but do

come from a location scale family so that

for any distribution

$$g(x) = \frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right) \text{ is the pdf.}$$

↑ rescaled version of a base pdf f

Then t & χ^2 are still pivotal, but they no longer

have the t_{n-1} and χ^2_{n-1} distributions respectively.

\Rightarrow We need different critical values.

Pivotal CIs for the mean μ

Let $t_{b,.025}$ and $t_{b,.975}$ satisfy

$$P(t_{b,.025} \leq t \leq t_{b,.975}) = 0.95.$$

Then $P(t_{b,.025} \leq \frac{\bar{X} - \mu}{s/\sqrt{n}} \leq t_{b,.975}) = 0.95$ Instead use t critical value for this formula. we gonna use critical value from bootstrap dist ~~test~~ for t

$$\Rightarrow P(t_{b,.025} \frac{s}{\sqrt{n}} \leq \bar{X} - \mu \leq t_{b,.975} \frac{s}{\sqrt{n}}) = 0.95$$

$$\Rightarrow P(-\bar{X} + t_{b,.025} \frac{s}{\sqrt{n}} \leq -\mu \leq -\bar{X} + t_{b,.975} \frac{s}{\sqrt{n}}) = 0.95$$

$$\Rightarrow P(\bar{X} - t_{b,.975} \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} - t_{b,.025} \frac{s}{\sqrt{n}}) = 0.95$$

$$\Rightarrow \left(\bar{X} - t_{b,.975} \frac{s}{\sqrt{n}}, \bar{X} - t_{b,.025} \frac{s}{\sqrt{n}} \right) \text{ is a}$$

95% CI for μ .

usually
negative

Pivotal CIs for the variance σ^2

By a similar argument, the 95% CI will be

$$\left(\frac{(n-1)s^2}{\chi^2_{b,.975}}, \frac{(n-1)s^2}{\chi^2_{b,.025}} \right)$$

97.5th percentile \rightarrow $\chi^2_{b,.975}$ \leftarrow 2.5th percentile $\chi^2_{b,.025}$

We estimate the needed percentiles for the distributions of t and χ^2 using the bootstrap.

Ex: Find bootstrap 95% pivotal confidence intervals for μ and σ^2 using the data from the earlier $n=10$ example.

See next page first.

Important Clarification

When we do the bootstrap to estimate the distribution

$$\text{for } t = \frac{\bar{X} - \mu}{s/\sqrt{n}} \text{ or } \chi^2 = \frac{(n-1)s^2}{\sigma^2}, \text{ the}$$

sample plays the role of the population, and the

bootstrap sample plays the role of the sample.

$$\underline{t}: \quad t_b = \frac{\bar{X}_b - \bar{X}}{s_b/\sqrt{n}} \leftarrow \text{from sample}$$

$$\underline{\chi^2}: \quad \chi_b^2 = \frac{(n-1)s_b^2}{s^2}$$

Comparison Example

Ex 60: Compare the performance of (normal-theory), (percentile bootstrap), and (bootstrap pivotal confidence intervals for the population variance).

In particular, how well do 95% confidence intervals maintain their level?

Use ① normal data and ② exponential data.
↳ 2nd part

Use sample sizes of 5, 10, 20, 40.

Q: What if we have uniform data?

Extending the Basic Bootstrap I

There are many different versions of the bootstrap for different kinds of data.

What can we do if we have bivariate or multivariate data?

Ex: Here's a sample of heights for randomly selected mother/daughter pairs. The units are inches.

$n = 10$ pairs

Mother	Daughter
70	67
69	64
65	62
64	64
66	69
65	70
64	65
66	66
60	63
70	74

Using the bootstrap, find

- (a) r , (b) an SE for r ,
- (c) an estimate of the bias of r ,

and

- (d) a bootstrap 95% CI for ρ , the population correlation.

Extending the Basic Bootstrap (II)

Ex: Ten subjects were randomly assigned to each of two treatments. Find

(a) A bootstrap SE for $\bar{X}_1 - \bar{X}_2$.

(b) A bootstrap 95% CI for $\mu_1 - \mu_2$

(c) A bootstrap 95% CI for σ_1^2 / σ_2^2 .

How would we do this?

Note: There are multiple possible approaches that one might use. A key thing for success of the bootstrap is to mimic the original sampling as much as possible.

Two-Sample Example

Treatment	Values				
1	9	12	12	14	17
	19	21	22	26	31
2	8	9	10	11	13
	13	19	21	22	24