

Resampling Methods: The Bootstrap

Rebecca C. Steorts, Duke University

STA 325, Chapter 5 ISL

Agenda

- ▶ Re-sampling Methods
- ▶ Cross Validation
- ▶ The Bootstrap

Re-sampling Methods

A re-sampling method involves repeatedly drawing samples from a **training data set** and refitting a model to obtain additional information about that model.

Example: Suppose we want to know the variability associated with a linear regression model.

1. Draw different samples from the training data
2. Fit a linear regression to **each sample**
3. Examine how the fits differ

Re-sampling Methods

In this module, we focus on cross-validation (CV) and the bootstrap.

- ▶ CV can be used to estimate the test error associated with a statistical learning method to evaluate its performance or to select a model's level of flexibility
 - ▶ The bootstrap most commonly measures the accuracy of a parameter estimate or of a given statistical learning method.
1. Model assessment: the process of evaluating a model's performance
 2. Model selection: the process of selecting the proper level of flexibility for a model

The Bootstrap

The bootstrap is a widely applicable and extremely powerful statistical tool that can be used to **quantify the uncertainty** associated with a given estimator or statistical learning method.

The Bootstrap

As a simple example, the bootstrap can be used to estimate the standard errors of the coefficients from a linear regression fit.

Of course, we can get these from packages, so this isn't particularly useful, but this is just one simple example of the bootstrap.

The power of the bootstrap lies in the fact that it can be easily applied to a wide range of statistical learning methods, including some for which a measure of variability is otherwise **difficult to obtain** and is **not automatically output** by statistical software.

Toy example: Investing

Suppose we wish to determine the best investment allocation under a simple model.

Later, we explore the use of the bootstrap to assess the variability associated with the regression coefficients in a linear model fit.

Toy example: Investing

Suppose that we wish to invest a fixed sum of money in two financial assets that yield returns of X and Y , where X and Y are random quantities.

We will invest a fraction α of our money in X , and will invest the remaining $1 - \alpha$ in Y .

Since there is variability associated with the returns on these two assets, we wish to choose α to minimize the total risk, or variance, of our investment.

Toy example: Investing

That is we want to minimize

$$\text{Var}(\alpha X + (1 - \alpha)Y).$$

One can show (exercise) that the value that minimizes the risk is given by

$$\alpha = \frac{\sigma_Y^2 - \sigma_{XY}}{\sigma_Y^2 + \sigma_X^2 - 2\sigma_{XY}}, \quad (1)$$

where $\sigma_X^2 = \text{Var}(X)$, $\sigma_Y^2 = \text{Var}(Y)$, and $\sigma_{XY} = \text{Cov}(X, Y)$.

Toy example: Investing

- ▶ In reality, $\sigma_X^2, \sigma_Y^2, \sigma_{XY}$ are unknown.
- ▶ We can compute estimates of these quantities $\hat{\sigma}_X^2, \hat{\sigma}_Y^2, \hat{\sigma}_{XY}$ using a data set that contains past measurements for X and Y .
- ▶ We can then estimate the value of α that minimizes the variance of our investment using

$$\hat{\alpha} = \frac{\hat{\sigma}_Y^2 - \hat{\sigma}_{XY}}{\hat{\sigma}_Y^2 + \hat{\sigma}_X^2 - 2\hat{\sigma}_{XY}}, \quad (2)$$

Toy example: Investing

- ▶ It is natural to wish to quantify the accuracy of our estimate of α .
- ▶ We can understand how this might work for simulated data but in general, we cannot apply this to real data since we cannot generate new samples from the original population (since its unknown).

The Bootstrap

The bootstrap approach allows us to use a computer to emulate the process of obtaining new sample sets, so that we can estimate the variability of $\hat{\alpha}$ without generating additional samples.

Rather than repeatedly obtaining independent data sets from the population, we instead obtain distinct data sets by repeatedly sampling observations from the original data set.

The Bootstrap

Suppose we have a simple dataset Z with n observations.

1. Randomly select n observations from the data set in order to produce a bootstrap data set, Z^{*1} .
 - ▶ The sampling is performed with replacement, which means that the same observation can occur more than once in the bootstrap data set.
2. We can use Z^{*1} to produce a new bootstrap estimate for α , which we call $\hat{\alpha}^{*1}$.

The Bootstrap (continued)

- ▶ This procedure is repeated B times for some large value of B , in order to produce B different bootstrap data sets,

$$Z^{*1}, Z^{*2}, \dots, Z^{*B}.$$

and B corresponding α estimates,

$$\hat{\alpha}^{*1}, \hat{\alpha}^{*2}, \dots, \hat{\alpha}^{*B}.$$

3. We can compute the standard error of these bootstrap estimates using the formula

$$SE_B(\hat{\alpha}) = \sqrt{\frac{1}{B-1} \sum_{r=1}^B (\hat{\alpha}^{*r} - \frac{1}{B} \sum_{r'=1}^B \hat{\alpha}^{*r'})^2}$$

4. This serves as an estimate of the standard error of $\hat{\alpha}$ estimated from the original data set.

The Bootstrap

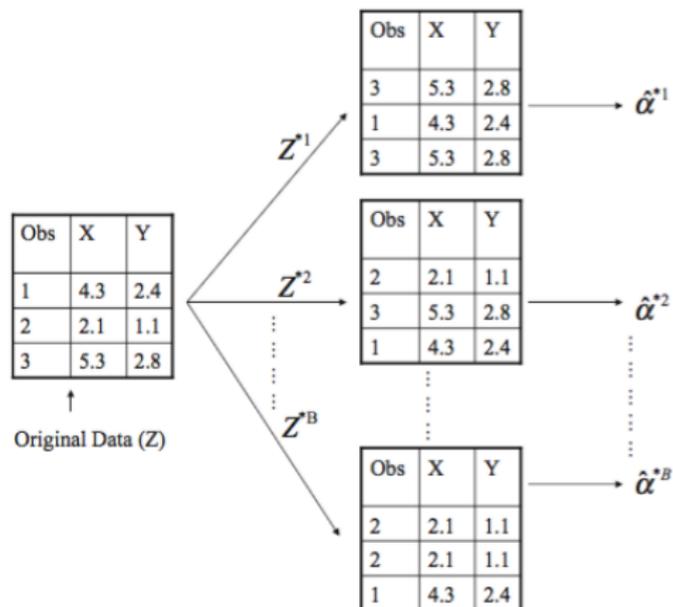


Figure 1: A graphical illustration of the bootstrap approach on a small sample containing $n = 3$ observations. Each bootstrap data set contains n observations, sampled with replacement from the original data set. Each bootstrap data set is used to obtain an estimate of α .

The Bootstrap in Practice

Performing a bootstrap analysis in R entails only two steps.

1. We must create a function that computes the statistic of interest.
2. We use the `boot()` function, which is part of the `boot` library, to perform the bootstrap by repeatedly sampling observations from the data set with replacement.

The Bootstrap on the Portfolio data set

The Portfolio data set in the ISLR package is the investment data set that motivated the bootstrap earlier.

To illustrate the use of the bootstrap on this data, we must

1. Create a function, `alpha.fn()`, which takes as input the (X, Y) data as well as a vector indicating which observations should be used to estimate α .
2. Then the function will output the estimate for α based on the selected observations.

The Bootstrap on the Portfolio data set

```
library(ISLR)
alpha.fn=function(data,index){
  X=data$X[index]
  Y=data$Y[index]
  return((var(Y)-cov(X,Y))/(var(X)+var(Y)-2*cov(X,Y)))
}
alpha.fn(Portfolio , 1:100)
```

```
## [1] 0.5758321
```

This function returns, or outputs, an estimate for α based on applying equation (5.7) to the observations indexed by the argument index.

The Bootstrap on the Portfolio data set

The next command uses the `sample()` function to randomly select 100 observations from the range 1 to 100, with replacement.

This is equivalent to constructing a new bootstrap data set and recomputing $\hat{\alpha}$ based on the new data set.

```
set.seed (1)
alpha.fn(Portfolio,sample(100,100,replace=T))
```

```
## [1] 0.5963833
```

Implementing the Bootstrap

We can implement a bootstrap analysis by performing this command many times, recording all of the corresponding estimates for α , and computing the resulting standard deviation.

The `boot()` function **automates** this approach.

Implementing the Bootstrap

Load the boot package in R.

```
library(boot)
```

Implementing the Bootstrap

```
# produce R=1000 bootstrap estimates  
# for alpha using boot()  
boot(Portfolio, alpha.fn, R=1000)
```

```
##  
## ORDINARY NONPARAMETRIC BOOTSTRAP  
##  
##  
## Call:  
## boot(data = Portfolio, statistic = alpha.fn, R = 1000)  
##  
##  
## Bootstrap Statistics :  
##      original      bias    std. error  
## t1* 0.5758321 -7.315422e-05 0.08861826
```

Bootstrap Summary Results

The final output shows that using the original data, $\hat{\alpha} = 0.5758$, and that the bootstrap estimate for $SE(\hat{\alpha})$ is 0.0886.

Bootstrap: Other Applied Examples with R

There is an excellent lab on applying the bootstrap to linear regression. Please work through this on your own. See ISL, page 195 – 197.