

towards
data science

Sign in

Get started

Follow

615K Followers

·

Editors' Picks

Features

Deep Dives

Gr

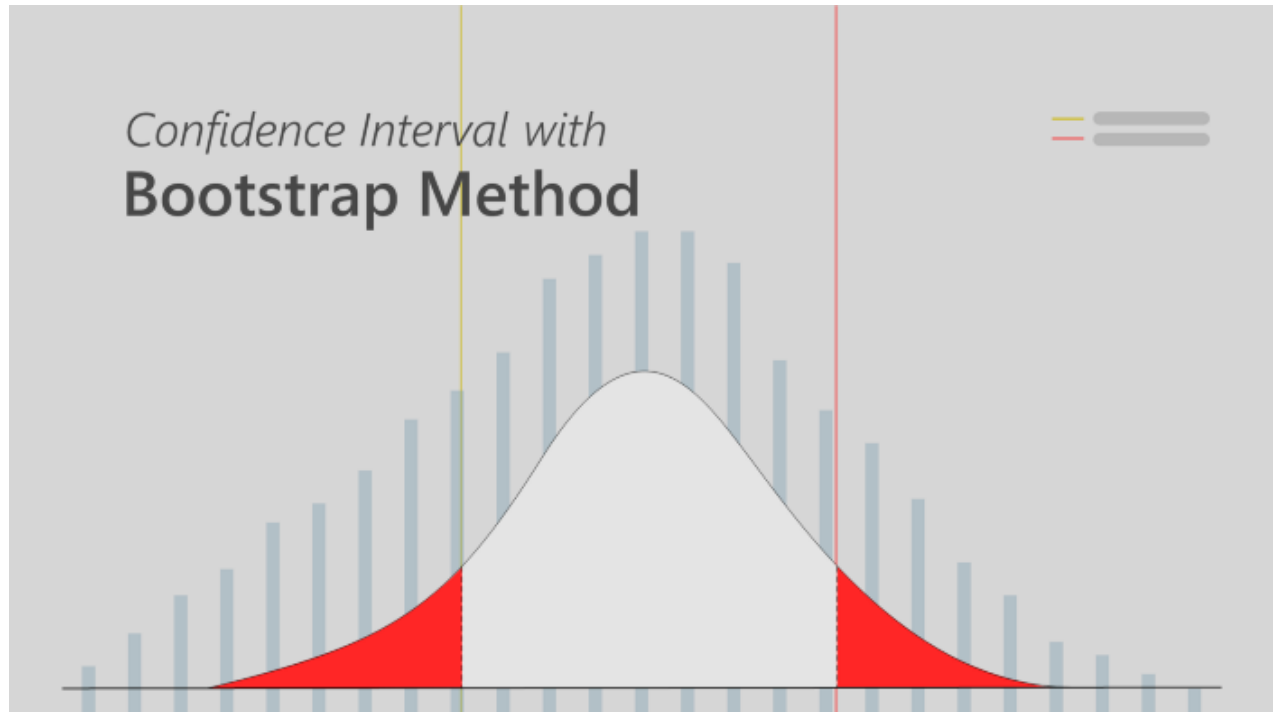
This is your **last** free member-only story this month. [Sign up for Medium and get an extra one](#)

Calculating Confidence Intervals with Bootstrapping

How can we calculate the confidence interval with bootstrapping?



Barış Hasdemir Jul 12, 2020 · 7 min read ★



Story banner, Image by author

Hi everyone,

In this article, I will attempt to explain how we can find a confidence interval by using Bootstrap Method. **Statistics** and **Python** knowledge are needed for better understanding.

Before diving into the method, let's remember some statistical concepts.

Variance: It is obtained by the sum of squared distances between a data point and the mean for each data point divided by the number of data points.

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$$

$$s^2 = \frac{\sum (x - \bar{x})^2}{n - 1}$$

Sample variance

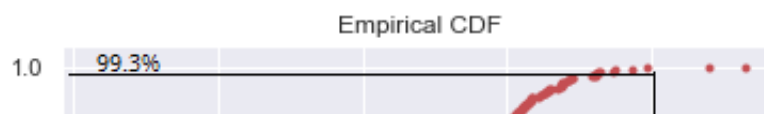
Standard Deviation: It is a measurement that shows us how our data points spread out from the mean. It is obtained by taking the square root of the variance

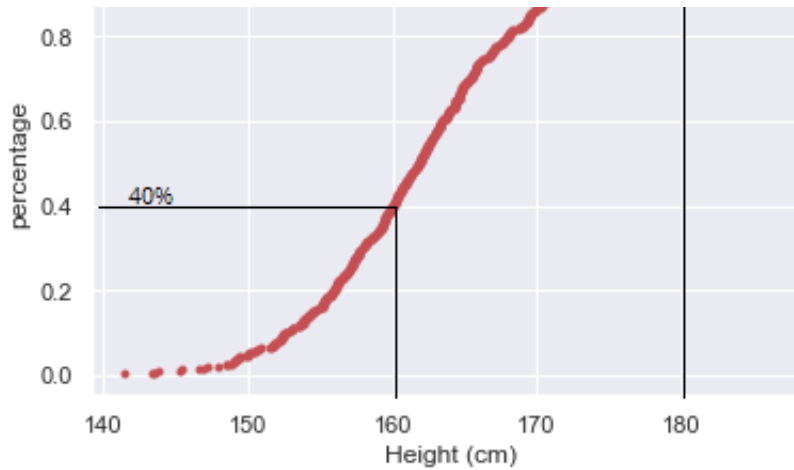
$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$$

Sample standard deviation

Cumulative Distribution Function: It can be used on any kind of variable X(discrete, continuous, etc.). It shows us the probability distribution of a variable. Therefore allowing us to interpret the probability of a value less than or equal to x from a given probability distribution

Empirical Cumulative Distribution Function: Also known as Empirical Distribution Function. The only difference between CDF and ECDF is, while the former shows us the hypothetical distribution of any given population, the latter is based on our observed data.





For example, how can we interpret the ECDF of the data shown on the chart above? We can say that 40% of heights are less than or equal to 160cm. Likewise, the percentage of people with heights of less than or equal to 180 cm is 99.3%

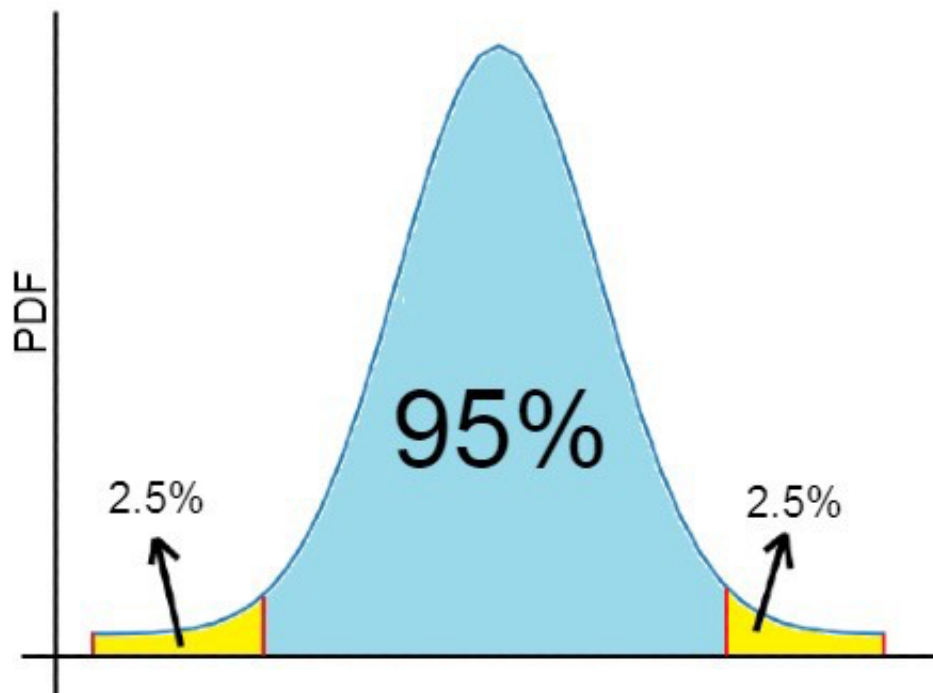
Probability Density Function: It shows us the distribution of continuous variables. The area under the curve gives us the probability so that the area must always be equal to 1

Normal Distribution: Also known as *Gaussian Distribution*. It is the most important probability distribution function in statistics which is bell-shaped and symmetric.

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Normal (Gaussian) Distribution

Confidence Interval: It is the range in which the values likely to exist in the population. It is estimated from the original sample and usually defined as 95% confidence but it may differ. You can consider the figure below which indicates a 95% confidence interval. The lower and upper limits of confidence interval defined by the values corresponding to the first and last 2.5th percentiles.



95% Confidence Interval, Image by author

What is Bootstrap Method?

Bootstrap Method is a resampling method that is commonly used in Data Science. It has been introduced by Bradley Efron in 1979. Mainly, it consists of the resampling our original sample with

replacement (*Bootstrap Sample*) and generating *Bootstrap replicates* by using *Summary Statistics*.

Confidence Interval of people heights

In this article, we are going to work with one of the datasets in Kaggle. It is *Weight-Height* data sets. It contains height (in inches) and weight (in pounds) information of 10.000 people separated by gender.

If you would like to see the whole code, you can find the *IPython notebook* via this [link](#).

We are going to use only heights of 500 randomly selected people and compute a 95% confidence interval by using Bootstrap Method

Let's start with importing the libraries that we will need.

```
1  # Import necessary libraries
2  import pandas as pd
3  import numpy as np
4  from scipy import stats
5  import seaborn as sns
6  import matplotlib.pyplot as plt
7
8  %matplotlib inline
9
10 # Enable Jupyter Notebook's intellisense
11 %config IPCompleter.greedy=True
```

```
## Bootstrapping a confidence interval, part 1
```

The first five rows of the DataFrame like following

```
1 # Import the iris dataset
2 data = pd.read_csv('weight-height.csv')
3
4 # Display first 5 rows
5 display(data.head())
```

Boostrapping2, ex hosted with a by GitHub

[view raw](#)

	Gender	Height	Weight
0	Male	73.847017	241.893563
1	Male	68.781904	162.310473
2	Male	74.110105	212.740856
3	Male	71.730978	220.042470
4	Male	69.881796	206.349801

Apparently, heights are in *inches*, let's convert heights from inches to *centimeters* and store in a new column *Height(cm)*.

```
1 # Convert inches to centimeters
2 data["Height(cm)"] = data["Height"]*2.54
3
4 # Get summary statistics of Heights in centimeters
5 display(data['Height(cm)'].describe())
```

Boostrapping2, ex hosted with a by GitHub

[view raw](#)

```
count    10000.000000
mean      168.573602
std        9.772721
min       137.828359
25%       161.304276
50%       168.447898
75%       175.702625
max       200.656806
Name: Height(cm), dtype: float64
```

As we can see above, the maximum and minimum height in the data set are 137.8 cm and 200.6 cm respectively.

We can use pandas.DataFrame's *sample* method to select 500 randomly selected heights. After that, we will print the summary statistics.

```
1 # Extract 500 random heights
2 heights = data['Height(cm)'].sample(500).reset_index(drop=True)
3
4 # Display Summary Statistics of heights in cm
5 display(heights.describe())
```

Bootstrapping? published with [ml by GitHub](#)

[view raw](#)

```
count      500.000000
mean       168.824438
std         9.824858
min        145.181405
25%        161.829489
50%        168.665300
75%        175.833922
max        198.363503
Name: Height(cm), dtype: float64
```

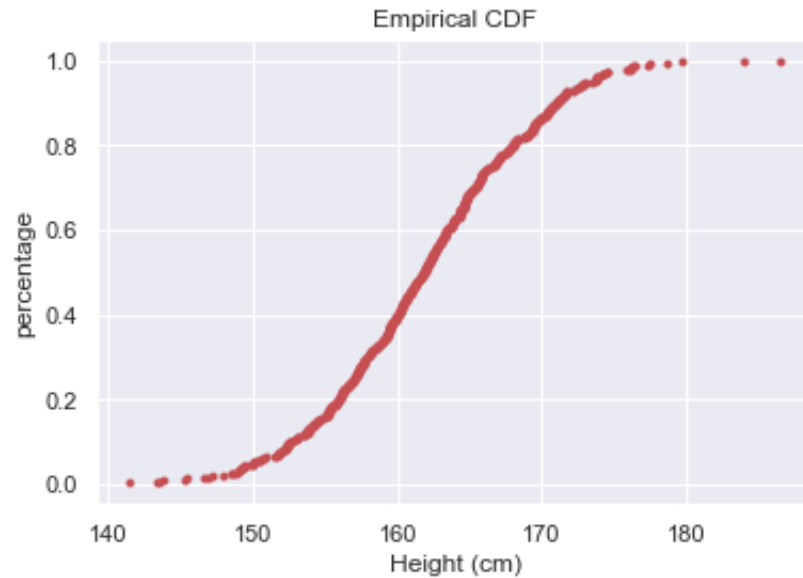

According to the output, our sample has 145 cm as minimum height and 198 cm as the maximum height.

Let's look at how ECDF and PDF look like?

```
1  # Create a function to get x, y for of ecdf
2  def get_ecdf(data):
3      """Returns x,y for ecdf"""
4      # Get lenght of the data into n
5      n = len(data)
6
7      # We need to sort the data
8      x = np.sort(data)
9
10     # the function will show us cumulative percentages of corresponding da
11     y = np.arange(1,n+1)/n
12
13     return x,y
14
15 # Create a function to plot ecdf
16 def plot_ecdf(data,labelx,labely,title,color):
17     """Plot ecdf"""
18     # Call get_ecdf function and assign the returning values
19     x, y = get_ecdf(data)
20
21     plt.plot(x,y,marker='.',linestyle='none',c=color)
22     plt.xlabel(labelx)
23     plt.ylabel(labely)
24     plt.title(title)
25
26 # Plotting Empirical CDF
27 plot_ecdf(heights,"Height (cm)","percentage","Empirical CDF","r")
28 plt.show()
```

ecdf.py hosted with ♥ by GitHub

[view raw](#)

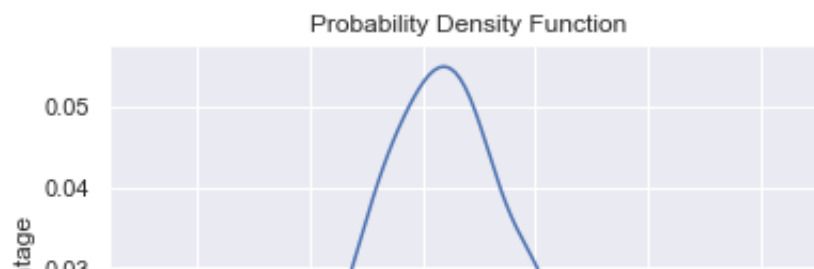


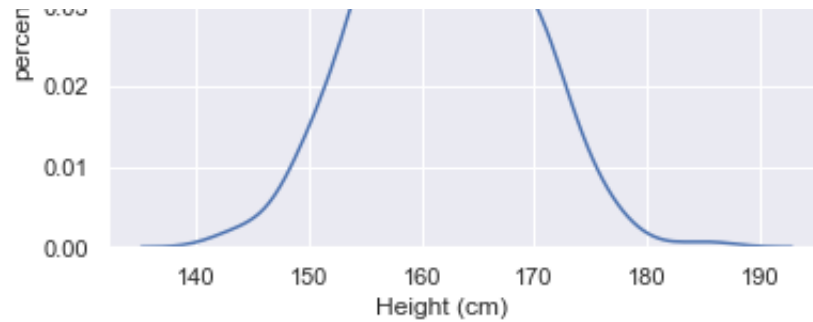
ECDF, Image by author

Empirical CDF demonstrates that 50% of people in our sample have 162 cm or less height.

What about PDF?

```
1 # Plotting PDF
2 sns.distplot(heights,hist=False)
3 plt.xlabel("Height(cm)")
4 plt.ylabel("PDF")
5 plt.title("Probability Density Function")
6 plt.show()
```



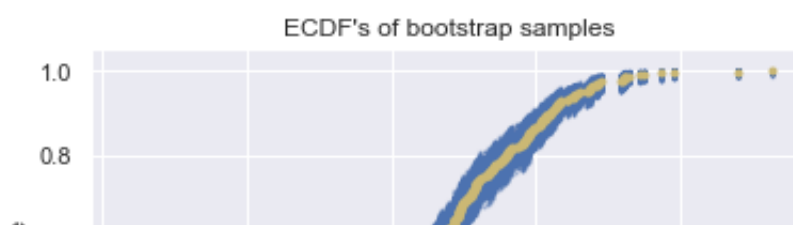


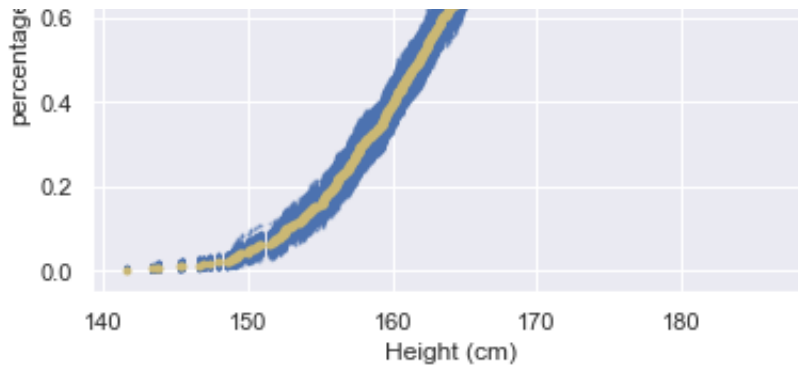
PDF, Image by author

PDF shows us the heights' distribution is too close to the normal distribution. Do not forget that the *area under the curve* gives us the probability.

Now, take a moment to think. We have only 500 observations in our sample, but there are billions of people in the world who we cannot measure their heights. Therefore, our sample does not give inference to the population. If we did the same measurements for different samples again and again, what would be the mean of heights?

For instance, assume that we did the same measurements with the same number of people (500) for 1000 times and plot the ECDF for each in a way that overlays the first observation's ECDF. It will look like the following.





ECDF, Image by author

As we can see above, we got different heights, but we can easily detect that the points are spreading in a specific range. That's the confidence interval that we want to learn

You may say that it is impossible to repeat the experiment so many times, you are not wrong. The exact reason why we use the Bootstrap Method. It helps us to simulate the same experiment thousands or even billions of times.

How?

In fact, the Bootstrap Method is quite straightforward and easy to understand. First, it generates bootstrap samples from our original sample by randomly choosing among the original sample. After that, it applies a summary statistics such as variation, standard deviation, mean, and so forth to get replicates. We will use 'mean' to generate our bootstrap replicates.

To understand the method, let's apply it to a small sample that

contains only 5 heights. We can generate our bootstrap samples like the following. Do not forget the fact that we can choose any observation more than once (resampling with replacement)

Barış Hasdemir

Data Engineer @trendyol



As we can see above we create 4 bootstrap samples and after that calculate their means. We will call these means our bootstrap replicates. Instead of 'mean' we could choose variance, standard deviation, median, or anything else.

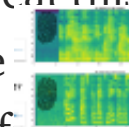
Related

Come back to our project. The next step, we are going to generate our bootstrap sample from our original sample to mean to get bootstrap replicate. We will repeat this process 15.000 times (drawing) in a for loop and store replicates in an array. To do this we can define a function like following

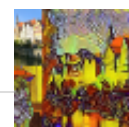
```
1 def draw_bs_replicates(data, func, size):
```



Andrew Chen on network effects and competing for growth in a pocketing market



A Speech Embedding Model for Speaker Recognition



Neutral Style Transfer
If you've ever pictured what ...

```

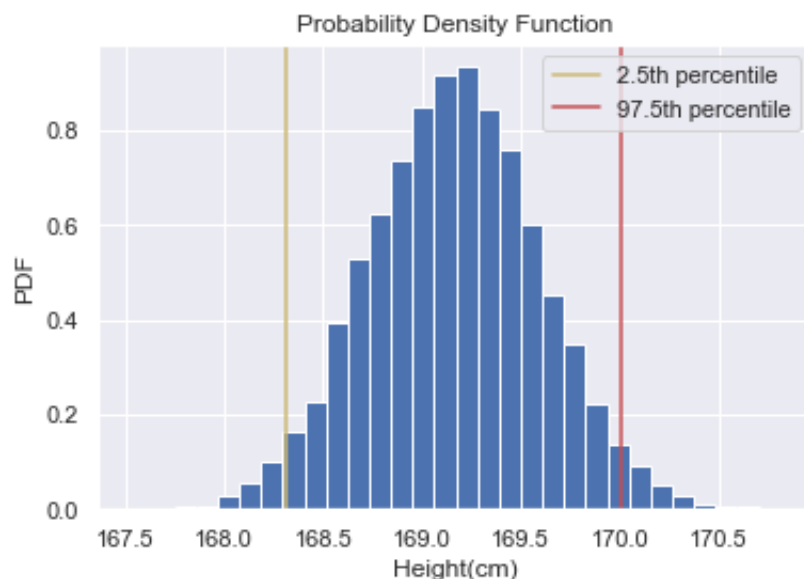
2 """creates a bootstrap sample, computes replicates and returns replicates
3 # Create an empty array to store replicates
4 bs_replicates = np.empty(size)
5
6 # Create bootstrap replicates as much as size
7 for i in range(size):
8     # Create a bootstrap sample
9     bs_sample = np.random.choice(data,size=len(data))
10    # Get bootstrap replicate and append to bs_replicates
11    bs_replicates[i] = func(bs_sample)
12
13    return bs_replicates

```

So, what are we going to do to calculate a 95% confidence interval?

After obtaining bootstrap replicates, the rest is so simple. As we know, our lower and upper limits are the values correspond to the 2.5th and 97.5th percentiles.

```
1 # Plot the PDF for bootstrap replicates as histogram
2 plt.hist(bs_replicates_heights,bins=30,normed=True)
3
4 # Showing the related percentiles
5 plt.axvline(x=np.percentile(bs_replicates_heights,[2.5]), ymin=0, ymax=1, color='yellow')
6 plt.axvline(x=np.percentile(bs_replicates_heights,[97.5]), ymin=0, ymax=1, color='red')
7
8 plt.xlabel("Height(cm)")
9 plt.ylabel("PDF")
10 plt.title("Probability Density Function")
11 plt.legend()
12 plt.show()
```



Percentiles, Image by author

We can find the boundaries with following simple Python code

```
1 # Get the corresponding values of 2.5th and 97.5th percentiles
2 conf_interval = np.percentile(bs_replicates_heights,[2.5,97.5])
3
4 # Print the interval
5 print("The confidence interval: ",conf_interval)
```

Bootstrapping6 published with [ml by GitHub](#)

[view raw](#)

```
The confidence interval: [167.77904671 169.51984003]
```

Our boundaries are found at 167.7 and 169.5. Therefore, we can say that if we do the same experiment with the whole population. The mean of heights will be between 167.7 cm and 169.5 cm with 95% of chance

Summary

Let's summarize what we did. We have randomly selected 500 heights and generated bootstrap samples. We calculated the 'mean' from those samples and got bootstrap replicates of means. Ultimately we calculated a 95% confidence interval.

I wish you good luck in your data journey :)

Reference

Introduction to Bootstrapping in Statistics with an Example — Statistics By Jim

Bootstrapping is a statistical procedure that resamples a single dataset to create many simulated samples....

statisticsbyjim.com

Statistical Thinking in Python (Part 2)

When doing statistical inference, we speak the language of probability. A probability distribution that...

www.datacamp.com

An Introduction to the Bootstrap Method

An exploration about bootstrap method, the motivation, and how it works

towardsdatascience.com

Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

Get this newsletter

[About](#) [Write](#) [Help](#) [Legal](#)