

Understanding Error Bars

A common error that learners run into with the week 3 assignment is looking at the *error bars of the data* rather than the *error bars of the means of the data*. These are very different, as the standard deviation of the means involves taking the square root of the number of samples.

This reading is intended to clarify the process required for assignment 3, with the demonstration based on the 1992 portion of the following data; we will create 1000 samples with a set random seed for reproducibility.

```
1  import pandas as pd
2  import numpy as np
3
4  df = pd.DataFrame([np.random.normal(32000,200000,3650),
5                    np.random.normal(43000,100000,3650),
6                    np.random.normal(43500,140000,3650),
7                    np.random.normal(48000,70000,3650)],
8                    index=[1992,1993,1994,1995])
9
10 # Let's do the random sampling 1000 times
11 np.random.seed(12345)
12 df_means = pd.DataFrame({'means': [np.random.normal(32000,200000,3650).mean() for i in range(1000)]})
13 df_means.head()
14
15 #means head output:
16 0    33312.107476
17 1    29723.719082
18 2    26276.149916
19 3    31267.288484
20 4    31121.673831
```

Using the 1000 samples of means, we will now compute the standard deviation.

```
1  df_means.std(axis=0)
2
3  #std output:
4  means    3414.816232
5  dtype: float64
```

This standard deviation is that of the means (also known as the standard error), and is the standard deviation used for the error bars. Note that this is not the standard deviation of the data.

The formula for calculating standard error is as follows ([see this Wikipedia article for more](#)^[7]):

$$s_{\bar{x}} = \frac{s}{\sqrt{n}}$$

Using the above formula, we can calculate the standard error as follows:

```
1  # data standard deviation: 200000
2  # sample size: 3650
3  import math
4  200000 / math.sqrt(3650)
5
6  #output:
7  3310.4235544094718
```