

Sampling Distributions

April 15, 2019

0.0.1 Sampling Distributions Introduction

In order to gain a bit more comfort with this idea of sampling distributions, let's do some practice in python.

Below is an array that represents the students we saw in the previous videos, where 1 represents the students that drink coffee, and 0 represents the students that do not drink coffee.

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
np.random.seed(42)

students = np.array([1,0,1,1,1,1,0,0,0,0,1,1,1,1,1,1,1,1,1,0])
```

1. Find the proportion of students who drink coffee in the above array. Store this value in a variable **p**.

```
In [2]: p = students.mean()
print(p)
```

```
0.714285714286
```

2. Use numpy's **random.choice** to simulate 5 draws from the `students` array. What is proportion of your sample drink coffee?

```
In [3]: np.random.choice(students,size =5, replace = True).mean()
```

```
Out[3]: 0.5999999999999999
```

3. Repeat the above to obtain 10,000 additional proportions, where each sample was of size 5. Store these in a variable called `sample_props`.

```
In [4]: sample_props = []
for _ in range(10000):
    sample = np.random.choice(students, 5, replace=True)
    sample_props.append(sample.mean())
```

4. What is the mean proportion of all 10,000 of these proportions? This is often called **the mean of the sampling distribution**.

```
In [5]: sample_props = np.array(sample_props)
        sample_props.mean()
```

```
Out[5]: 0.71399999999999997
```

5. What are the variance and standard deviation for the original 21 data values?

```
In [6]: student_var = students.var()
        student_std = students.std()
        print (student_var)
        print (student_std)
```

```
0.204081632653
```

```
0.451753951453
```

6. What are the variance and standard deviation for the 10,000 proportions you created?

```
In [7]: sample_var = sample_props.var()
        sample_var
```

```
Out[7]: 0.041763999999999996
```

```
In [8]: sample_std = sample_props.std()
        sample_std
```

```
Out[8]: 0.2043624231604235
```

7. Compute $p(1-p)$, which of your answers does this most closely match?

```
In [9]: p * (1-p)  # The variance of the original data
```

```
Out[9]: 0.20408163265306123
```

8. Compute $p(1-p)/n$, which of your answers does this most closely match?

```
In [10]: p*(1-p)/5 # The variance of the sample mean of size 5
```

```
Out[10]: 0.040816326530612249
```

9. Notice that your answer to 8. is commonly called the **variance of the sampling distribution**. If you were to change your first sample to be 20, what would this do for the variance of the sampling distribution? Simulate and calculate the new answers in 6. and 8. to check that the consistency you found before still holds.

```
In [11]: ##Simulate your 20 draws
        sample_props_20 = []
        for _ in range(10000):
            sample = np.random.choice(students, 20, replace=True)
            sample_props_20.append(sample.mean())
```

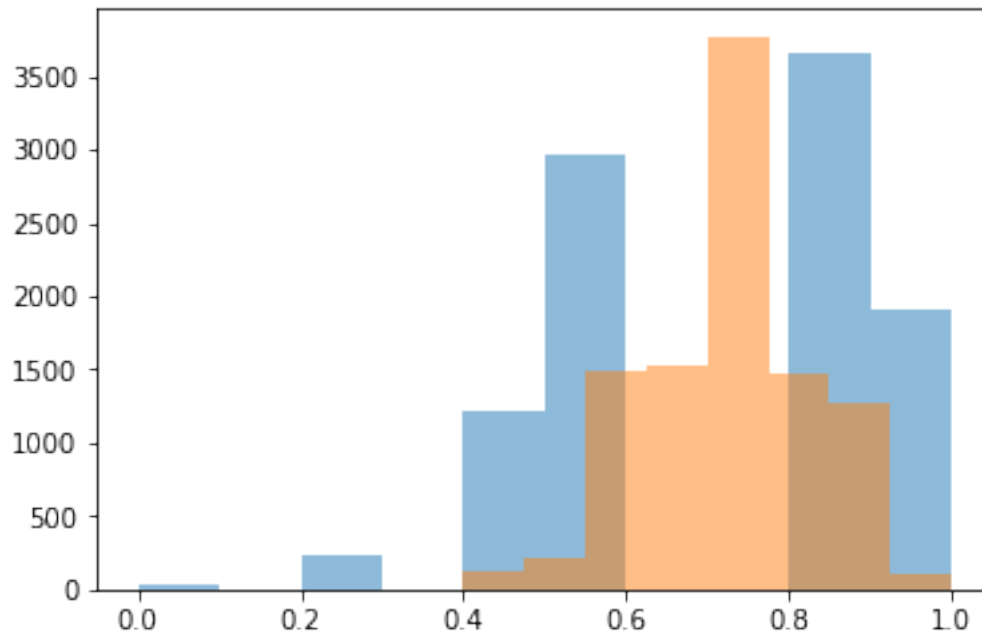
```
In [12]: ##Compare your variance values as computed in 6 and 8,
        ##but with your sample of 20 values
        print(p*(1-p)/20) # The theoretical variance
        print(np.array(sample_props_20).var()) # The simulated variance
```

0.0102040816327

0.010300994375

10. Finally, plot a histogram of the 10,000 draws from both the proportions with a sample size of 5 and the proportions with a sample size of 20. Each of these distributions is a sampling distribution. One is for the proportions of sample size 5 and the other a sampling distribution for proportions with sample size 20.

```
In [13]: plt.hist(sample_props, alpha=.5);
        plt.hist(np.array(sample_props_20), alpha=.5);
```



In []:

In []:

In []:

In []: