

CS 109A/STAT 121A/AC 209A/CSCI E-109A

Homework 0

Harvard University

Fall 2017

Instructors: Pavlos Protopapas, Kevin Rader, Rahul Dave, Margo Levine

This is a homework which you must turn in.

This homework has the following intentions:

1. To get you familiar with the jupyter/python environment (whether you are using your own install or jupyterhub)
2. You should easily understand these questions and what is being asked. If you struggle, this may not be the right class for you.
3. You should be able to understand the intent (if not the exact syntax) of the code and be able to look up google and provide code that is asked of you. If you cannot, this may not be the right class for you.

```
In [18]: # The line %... is a jupyter "magic" command, and is not part of the Python Language
# In this case we're just telling the plotting library to draw things on
# the notebook, instead of on a separate window.
%matplotlib inline
# See the "import ... as ..." constructs below? They're just aliasing the package
# That way we can call methods like plt.plot() instead of matplotlib.pyplot.plot()
import numpy as np
import matplotlib.pyplot as plt
```

Simulation of a coin throw

We don't have a coin right now. So let us **simulate** the process of throwing one on a computer. To do this we will use a form of the **random number generator** built into numpy. In particular, we will use the function `np.random.choice`, which will pick items with uniform probability from a list (thus if the list is of size 6, it will pick one of the six list items each time, with a probability 1/6).

```
In [19]: def throw_a_coin(N):
        return np.random.choice(['H','T'], size=N)

throws = throw_a_coin(40)
print("Throws",throws)
```

```
Throws ['H' 'T' 'T' 'H' 'T' 'T' 'H' 'T' 'T' 'H' 'T' 'H' 'H' 'H' 'T' 'T' 'H' 'T'
'H' 'T' 'T' 'H' 'H' 'T' 'H' 'T' 'H' 'H' 'T' 'T' 'T' 'H' 'H' 'T' 'T' 'H'
'H' 'H' 'T' 'T']
```

This next line gives you a True when the array element is a 'H' and False otherwise.

```
In [20]: throws == 'H'
```

```
Out[20]: array([ True, False, False,  True, False, False,  True, False, False,
        True, False,  True,  True,  True, False, False,  True, False,
        True, False, False,  True,  True, False,  True, False,  True,
        True, False, False, False,  True,  True, False, False,  True,
        True,  True, False, False], dtype=bool)
```

If you do a `np.sum` on the array of Trues and Falses, python will coerce the True to 1 and False to 0. Thus a sum will give you the number of heads

```
In [21]: np.sum(throws == 'H')
```

```
Out[21]: 19
```

```
In [22]: print("Number of Heads:", np.sum(throws == 'H'))
        print("p1 = Number of Heads/Total Throws:", np.sum(throws == 'H')/40.) # you can
```

```
Number of Heads: 19
p1 = Number of Heads/Total Throws: 0.475
```

Notice that you do not necessarily get 20 heads.

Now say that we run the entire process again, a second **replication** to obtain a second sample. Then we ask the same question: what is the fraction of heads we get this time? Lets call the odds of heads in sample 2, then, p_2 :

```
In [23]: throws = throw_a_coin(40)
        print("Throws:", throws)
        print("Number of Heads:", np.sum(throws == 'H'))
        print("p2 = Number of Heads/Total Throws:", np.sum(throws == 'H')/40.)
```

```
Throws: ['T' 'H' 'H' 'H' 'H' 'H' 'H' 'H' 'T' 'T' 'T' 'T' 'T' 'H' 'H' 'H' 'T' 'T'
        'T'
        'H' 'H' 'H' 'H' 'H' 'T' 'H' 'H' 'H' 'T' 'H' 'H' 'T' 'T' 'T' 'H' 'H' 'H'
        'H' 'H' 'H' 'T']
Number of Heads: 25
p2 = Number of Heads/Total Throws: 0.625
```

Q1. Show what happens as we choose a larger and larger set of trials

Do one replication for each size in the trials array below. Store the resultant probabilities in an array `probabilities`. Write a few lines on what you observe.

```
In [47]: trials = [10, 30, 50, 70, 100, 130, 170, 200, 500, 1000, 2000, 5000, 10000]
```

In [48]: probabilities = []

```
def rerun_trials(args):
    for arg in args:
        #probabilities.append(throw_a_coin(arg))
        throws = throw_a_coin(arg)
        print("Throws:", throws)
        print("Number of Heads:", np.sum(throws == 'H'))
        print("Number of Throws", arg)
        p = np.sum(throws == 'H')/arg
        print("p2 = Number of Heads/Total Throws:", p)
        print("\n")
        probabilities.append(p)

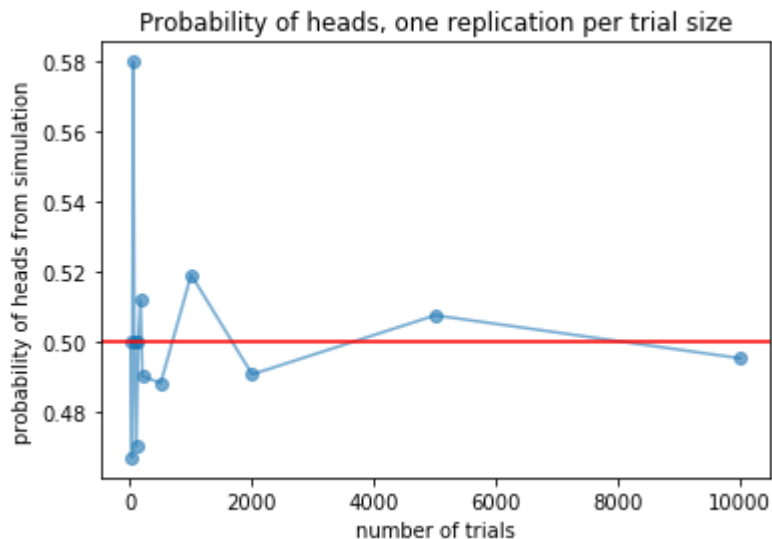
rerun_trials(trials)
print(probabilities)
```

Throws: ['H' 'T' 'T' 'H' 'H' 'T' 'H' 'T' 'T' 'H']
 Number of Heads: 5
 Number of Throws 10
 p2 = Number of Heads/Total Throws: 0.5

Throws: ['H' 'H' 'T' 'T' 'H' 'T' 'T' 'H' 'T' 'H' 'H' 'T' 'T' 'T' 'H' 'H' 'T'
 'T'
 'H' 'H' 'T' 'T' 'T' 'H' 'H' 'T' 'T' 'H' 'H' 'T']
 Number of Heads: 14
 Number of Throws 30
 p2 = Number of Heads/Total Throws: 0.4666666666667

Throws: ['T' 'T' 'H' 'H' 'T' 'T' 'T' 'H' 'H' 'H' 'T' 'H' 'T' 'T' 'T' 'H' 'H'
 'T'
 'H' 'H' 'H' 'H' 'H' 'H' 'T' 'H' 'H' 'T' 'H' 'H' 'H' 'H' 'T' 'H' 'H' 'H'
 'H' 'H' 'T' 'T' 'T' 'H' 'H' 'H' 'H' 'T' 'T' 'T' 'T' 'T']
 Number of Heads: 29
 Number of Throws 50

```
In [49]: plt.plot(trials, probabilities, 'o-', alpha=0.6);
plt.axhline(0.5, 0, 1, color='r');
plt.xlabel('number of trials');
plt.ylabel('probability of heads from simulation');
plt.title('Probability of heads, one replication per trial size');
```



What did you observe?

As the number of trials increase, the probability approaches .5

Multiple replications of the coin flips

Lets redo the experiment with coin flips that we started above. We'll establish some terminology at first. As notation we shall call the size of the trial of coin flips n . We'll call the result of each coin flip an observation, and a single replication (which is what we did above) a sample of observations. We will do M replications (or M "samples"), for which the variable in the function below is `number_of_samples` now, for each sample size n (`sample_size`).

Q2. Write a function to make M replications of N throws

Your job is to write a function `make_throws` which takes as arguments the `number_of_samples` (M) and the `sample_size` (n), and returns a list of probabilities of size M , with each probability coming from a different replication of size n . In each replication we do n coin tosses. We have provided a "spec" of the function below.

```

In [50]: """
Function
-----
make_throws

Generate a array of probabilities, each representing
the probability of finding heads in a sample of fair coins

Parameters
-----
number_of_samples : int
    The number of samples or replications
sample_size: int
    The size of each sample (we assume each sample has the same size)

Returns
-----
sample_probs : array
    Array of probabilities of H, one from each sample or replication

Example
-----
>>> make_throws(number_of_samples = 3, sample_size = 20)
[0.40000000000000002, 0.5, 0.59999999999999998]
"""

def make_throws(number_of_samples, sample_size):
    sample_probs = []
    while number_of_samples > 0:
        throws = throw_a_coin(sample_size)
        p = np.sum(throws == 'H')/sample_size
        sample_probs.append(p)
        number_of_samples -= 1
    return(sample_probs)

```

We show the mean over the observations, or sample mean, for a sample size of 10, with 20 replications. There are thus 20 means.

```
In [51]: make_throws(number_of_samples=20, sample_size=10)
```

```
Out[51]: [0.80000000000000004,
0.69999999999999996,
0.69999999999999996,
0.40000000000000002,
0.10000000000000001,
0.59999999999999998,
0.40000000000000002,
0.29999999999999999,
0.59999999999999998,
0.40000000000000002,
0.40000000000000002,
0.5,
0.69999999999999996,
0.40000000000000002,
0.69999999999999996,
0.40000000000000002,
0.5,
0.5,
0.40000000000000002,
0.69999999999999996]
```

Q3. What happens to the mean and standard deviation of the sample means as you increase the sample size

Using the sample sizes from the `sample_sizes` array below, compute a set of `sample_means` for each sample size, and for 200 replications. Calculate the mean and standard deviation for each sample size. Store this in arrays `mean_of_sample_means` and `std_dev_of_sample_means`. The standard deviation of the sampling means is called the "standard error". Explain what you see about this "mean of sampling means".

```
In [52]: sample_sizes = np.arange(1,1001,1)
```

```
In [87]: mean_of_sample_means = np.array([])
std_dev_of_sample_means = np.array([])
```

```
In [88]: # Iterate over sample sizes, calculating mean and standard deviation for each sam
for sample_size in sample_sizes:
    mean_of_sample_means = np.append(mean_of_sample_means, np.mean(make_throws(20
    std_dev_of_sample_means = np.append(std_dev_of_sample_means, np.std(make_thro
```

In [106]:

```

0.02260103, 0.02191433, 0.02210905, 0.02151248, 0.0208674 ,
0.02255985, 0.0219637 , 0.02201661, 0.02330411, 0.02040363,
0.02003377, 0.01988042, 0.02267737, 0.01926981, 0.01901649,
0.02004147, 0.02245176, 0.02196064, 0.02222295, 0.02036453,
0.02260996, 0.02290842, 0.0213811 , 0.02116235, 0.02159372,
0.02076352, 0.02009808, 0.01980762, 0.02208598, 0.02280365,

0.02053099, 0.02177152, 0.02212294, 0.02000941, 0.02293874,
0.02021125, 0.02128467, 0.01947438, 0.01931859, 0.01976401,
0.02043703, 0.02044667, 0.02076289, 0.01989476, 0.02148296,
0.02148327, 0.02004558, 0.02028528, 0.0195132 , 0.02178027,
0.01890233, 0.01997278, 0.02116303, 0.0202127 , 0.02165682,
0.02114852, 0.01984064, 0.02326219, 0.02062906, 0.02131379,
0.02178343, 0.02066722, 0.02009705, 0.01980306, 0.01970677,
0.02035382, 0.02096977, 0.01839424, 0.02103204, 0.02169213,
0.01965713, 0.01994509, 0.01961656, 0.02043239, 0.02168133,
0.0192964 , 0.02059525, 0.01937599, 0.01967999, 0.01972545,
0.01847234, 0.01965952, 0.01992601, 0.0186676 , 0.02106477,
0.0208111 , 0.02002697, 0.01983759, 0.02048176, 0.0185262 ,
0.0197804 , 0.0205755 , 0.02033848, 0.01964562, 0.01898635,
0.01867305, 0.02046652, 0.02046730, 0.02000386, 0.02105344

```

In [90]: *# mean and std of 200 means from 200 replications, each of size 10*
trials[0], mean_of_sample_means[0], std_dev_of_sample_means[0]

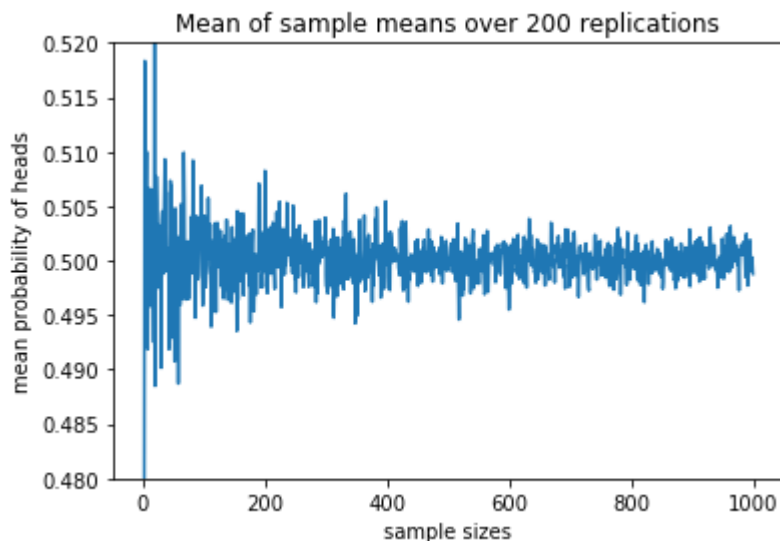
Out[90]: (10, 0.44, 0.4997749493522059)

In [120]:

```

plt.plot(sample_sizes, mean_of_sample_means);
plt.ylim([0.480,0.520]);
plt.xlabel("sample sizes")
plt.ylabel("mean probability of heads")
plt.title("Mean of sample means over 200 replications");

```



Explain what you see about this "mean of sampling means".

The mean of the sample means approaches .5 as the sample size increases, but there is still variance around .5

Q4. What distribution do the sampling means follow?

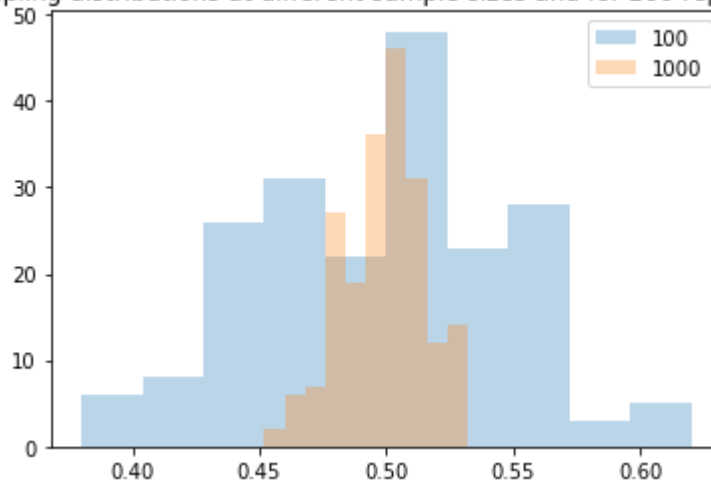
Store in variables `sampling_means_at_size_100` and `sampling_means_at_size_1000` the set of sampling means at sample sizes of 100 and 1000 respectively, still with 200 replications. We will plot in a histogram below these distributions. What type of distributions are these, roughly? How do these distributions vary with sample size?

```
In [121]: sampling_means_at_size_100 = []
sampling_means_at_size_1000 = []
sampling_means_at_size_100 = make_throws(200, 100)
sampling_means_at_size_1000 = make_throws(200, 1000)
```

```
In [ ]:
```

```
In [122]: plt.hist(sampling_means_at_size_100, alpha=0.3, label="100", bins=10)
plt.hist(sampling_means_at_size_1000, alpha=0.3, label="1000", bins=10)
plt.legend();
plt.title("Sampling distributions at different sample sizes and for 200 replications")
```

Sampling distributions at different sample sizes and for 200 replications



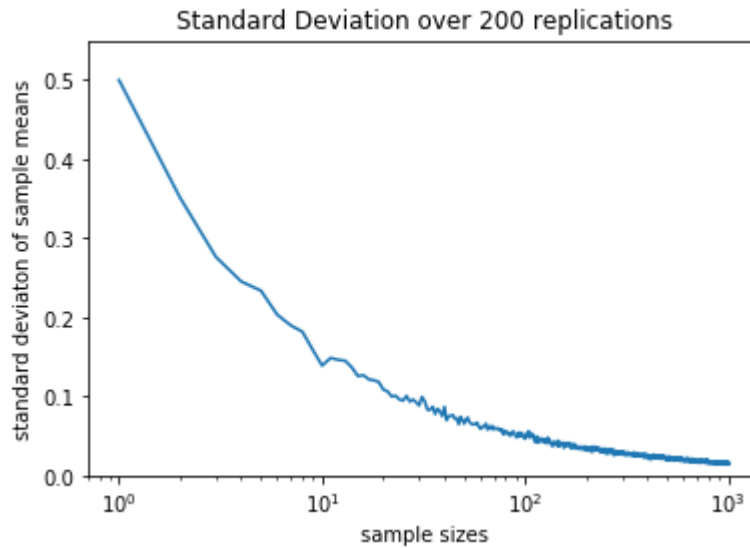
What type of distributions are these, roughly? How do these distributions vary with sample size?

As sample sizes increase, the distribution approaches a normal distribution

Q5. How does the standard error of the sample mean vary with sample size? Create a plot to illustrate how it varies over various sample sizes.

Hint: you might want to take logarithms for one of your axes


```
In [124]: plt.plot(sample_sizes, std_dev_of_sample_means);  
plt.ylim([0.000,0.550]);  
plt.xlabel("sample sizes")  
plt.ylabel("standard deviaton of sample means")  
plt.title("Standard Deviation over 200 replications");  
plt.xscale('log');
```



How does the standard error of the sample mean vary with sample size?

As the sample size increases, the standard error decreases from the mean towards 0 and levels off.