



CS109A Introduction to Data Science

Homework 0

Harvard University

Summer 2018

Instructors: Pavlos Protopapas and Kevin Rader

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.
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 on Google and provide code that is asked of you. If you cannot, this may not be the right class for you.

```
In [1]: ## RUN THIS CELL TO GET THE RIGHT FORMATTING
from IPython.core.display import HTML
def css_styling():
    styles = open("cs109.css", "r").read()
    return HTML(styles)
css_styling()
```

Out[1]:

Basic Math and Probability/Statistics Calculations

We'll start you off with some basic math and statistics problems questions to make sure you have the appropriate background to be comfortable with concepts that will come up in CS 109a.

Question 1: Mathiness is What Brings Us Together Today

Matrix Operations

Complete the following matrix operations. **Note: Do not do this numerically and show your work as a markdown/latex notebook cell**

1.1. Let $A = \begin{pmatrix} 3 & 4 & 2 \\ 5 & 6 & 4 \\ 4 & 3 & 4 \end{pmatrix}$ and $B = \begin{pmatrix} 1 & 4 & 2 \\ 1 & 9 & 3 \\ 2 & 3 & 3 \end{pmatrix}$.

Compute $A \cdot B$.

1.2. Let $A = \begin{pmatrix} 0 & 12 & 8 \\ 1 & 15 & 0 \\ 0 & 6 & 3 \end{pmatrix}$.

Compute A^{-1} .

Calculus and Probability

Complete the following (show your work as a markdown/latex notebook cell)

1.3. From Wikipedia:

In mathematical optimization, statistics, econometrics, decision theory, machine learning and computational neuroscience, a loss function or cost function is a function that maps an event or values of one or more variables onto a real number intuitively representing some "cost" associated with the event. An optimization problem seeks to minimize a loss function.

We've generated a cost function on parameters $x, y \in \mathbb{R}$ $L(x, y) = 3x^2y - y^3 - 3x^2 - 3y^2 + 2$. Find the critical points (optima) of $L(x, y)$.

1.4. A central aspect of call center operations is the per minute statistics of caller demographics. Because of the massive call volumes call centers achieve, these per minute statistics can often take on well-known distributions. In the CS109 Homework Helpdesk, X and Y are discrete random variables with X measuring the number of female callers per minute and Y the total number of callers per minute. We've determined historically the joint pmf of (X, Y) and found it to be

$$p_{X,Y}(x,y) = e^{-4} \frac{2^y}{x!(y-x)!}$$

where $y \in \mathbb{N}, x \in [0, y]$ (That is to say the total number of callers in a minute is a non-negative integer and the number of female callers naturally assumes a value between 0 and the total number of callers inclusive). Find the mean and variance of the marginal distribution of X .

*Hints: *

1. $x \in [0, y] \Rightarrow x < y$.
2. You may find the change of variable $z = y - x$ helpful.
3. Recall:

$$\sum_{z=0}^{\infty} \frac{2^z}{z!} = e^2$$

Basic Statistics

Complete the following: you can perform the calculations by hand (show your work) or using software (include the code and output...screenshots are fine if it is from another platform).

1.5. 37 of the 76 female CS concentrators have taken Data Science 1 (DS1) while 50 of the 133 male concentrators haven taken DS1. Perform a statistical test to determine if interest in Data Science (by taking DS1) is related to sex. Be sure to state your conclusion.

Answers

1.1

We can rewrite A and B in terms of row and column vectors respectively. Let

$$A = \begin{bmatrix} r_1 \\ r_2 \\ r_3 \end{bmatrix}$$

and

$$B = \begin{bmatrix} c_1 & c_2 & c_3 \end{bmatrix}.$$

Then

$$\begin{aligned} A \cdot B &= \begin{bmatrix} r_1 \cdot c_1 & r_1 \cdot c_2 & r_1 \cdot c_3 \\ r_2 \cdot c_1 & r_2 \cdot c_2 & r_2 \cdot c_3 \\ r_3 \cdot c_1 & r_3 \cdot c_2 & r_3 \cdot c_3 \end{bmatrix} \\ &= \begin{bmatrix} 11 & 54 & 24 \\ 19 & 86 & 40 \\ 15 & 55 & 29 \end{bmatrix} \end{aligned}$$

1.2

First, form an augmented matrix with A and the 3x3 identity matrix

$$A = \begin{bmatrix} 0 & 12 & 8 & 1 & 0 & 0 \\ 1 & 15 & 0 & 0 & 1 & 0 \\ 0 & 6 & 3 & 0 & 0 & 1 \end{bmatrix}$$

The reduced row echelon form of this matrix is

$$\begin{bmatrix} 1 & 0 & 0 & 15/4 & 1 & -10 \\ 0 & 1 & 0 & -1/4 & 0 & 2/3 \\ 0 & 0 & 1 & 1/2 & 0 & -1 \end{bmatrix}$$

This means that

$$A^{-1} = \begin{bmatrix} 5/4 & 1 & -10 \\ -1/4 & 0 & 2/3 \\ 1/2 & 0 & -1 \end{bmatrix}$$

1.3

To find the critical points we take the partial derivatives of L

$$L_x = 6x(y - 1)$$

$$L_y = 3x^2 - 3y^2 - 6y$$

L_x is equal to 0 when $x = 0$ or $y = 1$. Setting $x = 0$,

$$\begin{aligned} L_y|_{x=0} &= 3y^2 - 6y \\ &= 3y(y - 2) \end{aligned}$$

Hence, (0, 0) and (0, 2) are stable points of L.

Alternatively, if we set $y = 1$,

$$L_y|_{y=1} = 3x^2 - 9$$

This means that we also have stable points at $(\sqrt{3}, 1)$ and $(-\sqrt{3}, 1)$

To determine which any of these stable points are saddle points, we use the 2-dimensional version of the 2nd partial derivative test.

Let H be the Hessian of L. We find the determinant of H.

$$\begin{aligned} |H| &= \begin{vmatrix} 6(y-1) & 6x \\ 6x & 6(-y-1) \end{vmatrix} \\ &= 36(x^2 + y^2 - 1) \end{aligned}$$

The determinant of the Hessian evaluated at $(\sqrt{3}, 1)$, $(-\sqrt{3}, 1)$ (0, 2) is positive, meaning that these stable points are optima. The determinant of the Hessian evaluated at (0, 0) is negative, meaning that this stable point is a saddle point.

The optima of L are $(\sqrt{3}, 1)$, $(-\sqrt{3}, 1)$ $(0, 2)$

1.4

The marginal probability density function for X can be found by summing over all values of y . Since $x \leq y$,

$$p_X(x) = \sum_{y=x}^{\infty} e^{-4} \frac{2^y}{x!(y-x)!}$$

We can write the expectation as,

$$\begin{aligned} E(X) &= \sum_{x=0}^{\infty} x p_X(x) \\ &= \sum_{x=0}^{\infty} x \sum_{y=x}^{\infty} e^{-4} \frac{2^y}{x!(y-x)!} \\ &= e^{-4} \sum_{x=0}^{\infty} \sum_{y=x}^{\infty} x \frac{2^x 2^{y-x}}{x!(y-x)!} \end{aligned}$$

If we let $z = y - x$, the expected value of X is

$$\begin{aligned} &= e^{-4} \sum_{x=0}^{\infty} x \frac{2^x}{x!} \sum_{z=0}^{\infty} \frac{2^z}{z!} \\ &= e^{-2} \sum_{x=0}^{\infty} x \frac{2^x}{x!} \\ &= e^{-2} \sum_{x=1}^{\infty} x \frac{2^x}{x!} \\ &= 2e^{-2} \sum_{x=1}^{\infty} \frac{2^{x-1}}{(x-1)!} \\ &= 2 \end{aligned}$$

By the same line of reasoning,

$$\begin{aligned}
E(X^2) &= \sum_{x=0}^{\infty} x^2 p_X(x) \\
&= \sum_{x=0}^{\infty} x^2 \sum_{y=x}^{\infty} e^{-4} \frac{2^y}{x!(y-x)!} \\
&= e^{-4} \sum_{x=0}^{\infty} \sum_{y=x}^{\infty} x^2 \frac{2^x 2^{y-x}}{x!(y-x)!} \\
&= e^{-4} \sum_{x=0}^{\infty} x^2 \frac{2^x}{x!} \sum_{z=0}^{\infty} \frac{2^z}{z!} \\
&= e^{-2} \sum_{x=0}^{\infty} x^2 \frac{2^x}{x!} \\
&= e^{-2} \left(\sum_{x=2}^{\infty} x^2 \frac{2^x}{x!} + 2 \right) \\
&= e^{-2} \left(4 \sum_{x=2}^{\infty} \frac{2^{x-2}}{(x-2)!} + 2 \right) \\
&= 4 + 2e^{-2}
\end{aligned}$$

Hence,

$$\begin{aligned}
\text{Var}(X) &= E(X^2) - E(X)^2 \\
&= 2e^{-2}
\end{aligned}$$

1.5

Let p_f represent the proportion of female students who are interested in data science and p_m represent the proportion of male students who are interested in data science. Our null hypothesis is $p_f - p_m = 0$.

We assume that $\hat{p}_f - \hat{p}_m$ is approximately normally distributed with mean $p_f - p_m$ and variance

$\frac{p_f(1-p_f)}{n_f} + \frac{p_m(1-p_m)}{n_m}$, which under the null hypothesis is $p(1-p) \left(\frac{1}{n_f} + \frac{1}{n_m} \right)$. We can standardize to get

$$\begin{aligned}
Z &= \frac{(\hat{p}_f - \hat{p}_m) - (p_f - p_m)}{\sqrt{p(1-p) \left(\frac{1}{n_f} + \frac{1}{n_m} \right)}} \\
&= \frac{(\hat{p}_f - \hat{p}_m)}{\sqrt{p(1-p) \left(\frac{1}{n_f} + \frac{1}{n_m} \right)}}
\end{aligned}$$

To estimate p , we define $\hat{p} = \frac{d_f + d_m}{n_f + n_m}$, where d_f is the number of female students in DS1 and d_m is the number of male students in DS1.

Thus, the test statistic is

$$Z = \frac{(\hat{p}_f - \hat{p}_m)}{\sqrt{\hat{p}(1 - \hat{p})\left(\frac{1}{n_f} + \frac{1}{n_m}\right)}}$$

Since $\hat{p}_f = 37/76 = 0.487$, $\hat{p}_m = 50/133 = 0.376$ and $\hat{p} = \frac{37+50}{76+133} = 0.416$, we can calculate the test statistic to be

$$Z = 1.566$$

Under a confidence level of 95%, i.e. $Z \in [-1.96, 1.96]$, we fail to reject our null hypothesis and conclude that gender is not related to interest in Data Science.

```
In [2]: ## RUN THIS CELL
# The line %... is a jupyter "magic" command, and is not part of the Python
# 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 packages
# That way we can call methods like plt.plot() instead of matplotlib.pyplot.
import numpy as np
import scipy as sp
import scipy.stats
import matplotlib.pyplot as plt
```

Simulation of a Coin Throw

We'd like to do some experiments with coin flips, but we don't have a physical coin at the moment. So let's **simulate** the process of flipping a coin 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 picks items with uniform probability from a list. If we provide it a list `['H', 'T']`, it will pick one of the two items in the list. We can also ask it to do this multiple times by specifying the parameter `size`.

```
In [3]: def throw_a_coin(n_trials):
        return np.random.choice(['H', 'T'], size=n_trials)
```

`np.sum` is a function that returns the sum of items in an iterable (i.e. a list or an array). Because python coerces `True` to 1 and `False` to 0, the effect of calling `np.sum` on the array of `True`s and `False`s will be to return the number of `True`s in the array (which can then effectively count the number of heads).

Question 2: The 12 Labors of Bernoullis

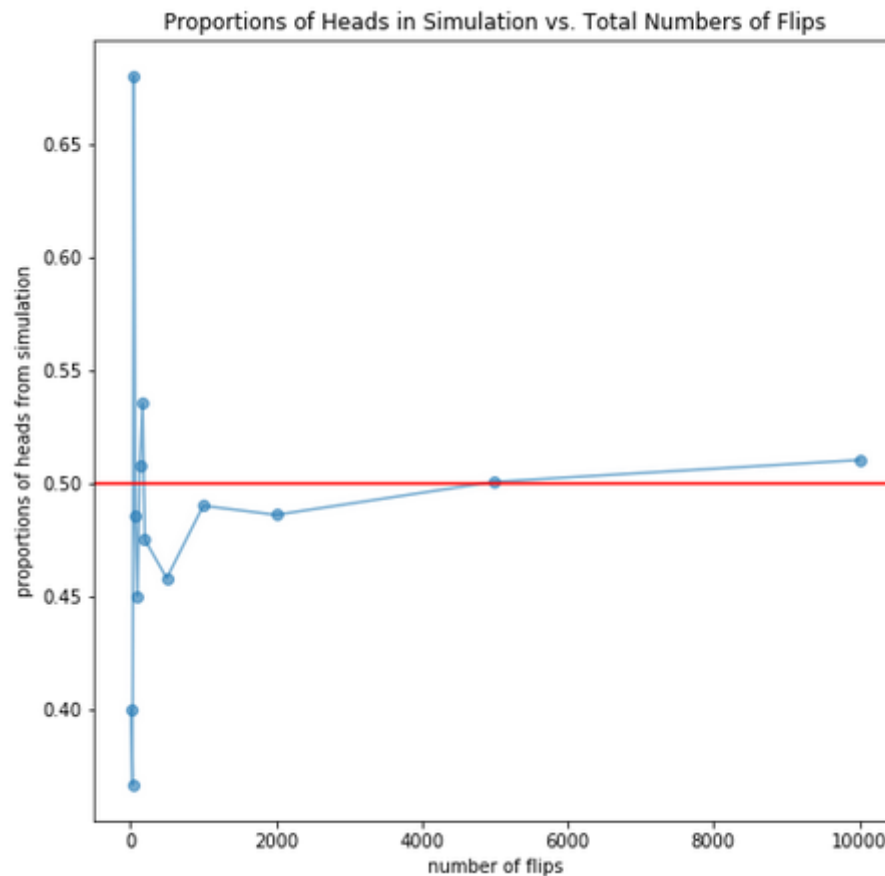
Now that we know how to run our coin flip experiment, we're interested in knowing what happens as we choose larger and larger number of coin flips.

2.1. Run one experiment of flipping a coin 40 times storing the resulting sample in the variable `throws1`. What's the total proportion of heads?

2.2. Replicate the experiment in 2.1 storing the resulting sample in the variable `throws2`. What's the proportion of heads? How does this result compare to that you obtained in question 2.1?

2.3. Write a function called `run_trials` that takes as input a list, called `n_flips`, of integers representing different values for the number of coin flips in a trial. For each element in the input list, `run_trials` should run the coin flip experiment with that number of flips and calculate the proportion of heads. The output of `run_trials` should be the list of calculated proportions. Store the output of calling `run_trials` in a list called `proportions`.

2.4. Using the results in 2.3, reproduce the plot below.



2.5. What's the appropriate observation about the result of running the coin flip experiment with larger and larger numbers of coin flips? Choose the appropriate one from the choices below and explain why.

A. Regardless of sample size the probability of in our experiment of observing heads is 0.5 so the proportion of heads observed in the coin-flip experiments will always be 0.5.

B. The proportions **fluctuate** about their long-run value of 0.5 (what you might expect if you tossed the coin an infinite amount of times), in accordance with the notion of a fair coin (which we encoded in our simulation by having `np.random.choice` choose between two possibilities with equal probability), with the fluctuations seeming to become much smaller as the number of trials increases.

C. The proportions **fluctuate** about their long-run value of 0.5 (what you might expect if you tossed the coin an infinite amount of times), in accordance with the notion of a fair coin (which we encoded in our simulation by having `np.random.choice` choose between two possibilities with equal probability), with the fluctuations constant regardless of the number of trials.

Answers

2.1

```
In [4]: def find_prop_heads(throws, n_trials):
        num_of_heads = np.sum(throws == 'H')
        prop_heads = float(num_of_heads)/float(n_trials)
        return prop_heads

n_trials = 40
throws1 = throw_a_coin(n_trials)
print('The total proportion of heads is %f' % find_prop_heads(throws1, n_trials))

The total proportion of heads is 0.475000
```

2.2

```
In [5]: n_trials = 40
        throws2 = throw_a_coin(n_trials)
        print('The total proportion of heads is %f' % find_prop_heads(throws2, n_trials))

The total proportion of heads is 0.550000
```

2.3

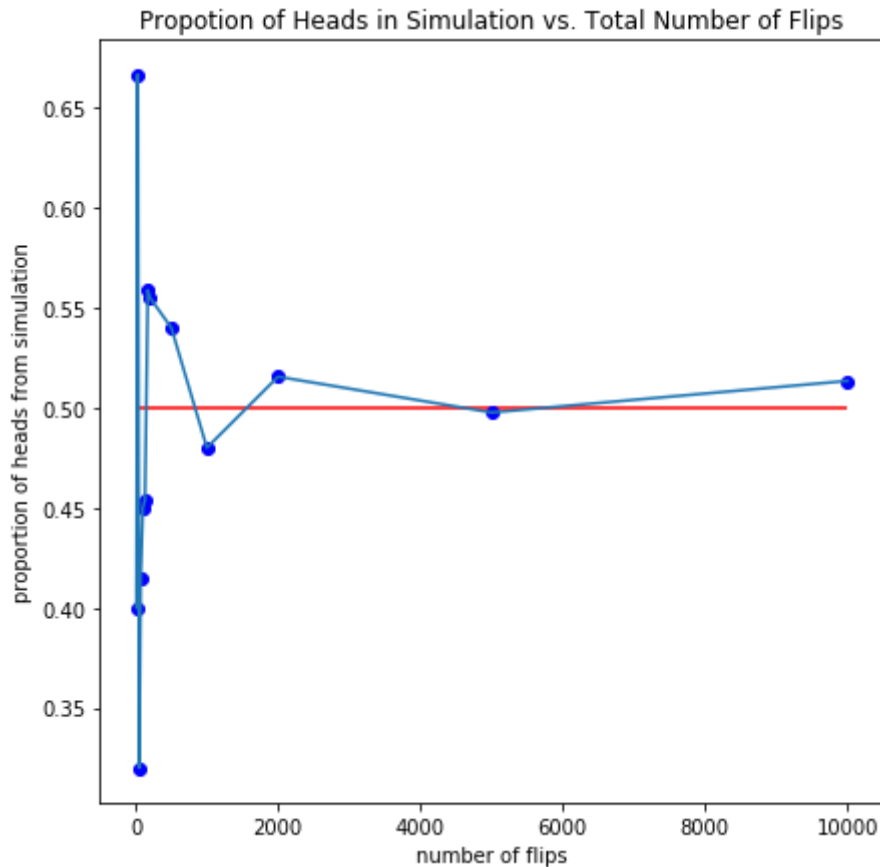
```
In [6]: n_flips = [10, 30, 50, 70, 100, 130, 170, 200, 500, 1000, 2000, 5000, 10000]
```

```
In [7]: def run_trials(n_flips_list):  
        return [find_prop_heads(throw_a_coin(n_trials), n_trials) for n_trials in  
  
proportions = run_trials(n_flips)  
proportions
```

```
Out[7]: [0.4,  
         0.6666666666666666,  
         0.32,  
         0.4142857142857143,  
         0.45,  
         0.45384615384615384,  
         0.5588235294117647,  
         0.555,  
         0.54,  
         0.48,  
         0.5155,  
         0.4976,  
         0.5134]
```

2.4

```
In [8]: plt.figure(figsize=(7, 7))
plt.plot(n_flips, proportions, 'bo', n_flips, proportions)
plt.ylabel('proportion of heads from simulation')
plt.xlabel('number of flips')
plt.title('Propotion of Heads in Simulation vs. Total Number of Flips')
plt.hlines(0.5, 0, 10000, color = 'r')
plt.show()
```



2.5

What's the appropriate observation about the result of applying the coin flip experiment to larger and larger numbers of coin flips? Choose the appropriate one.

B. The proportions **fluctuate** about their long-run value of 0.5 (what you might expect if you tossed the coin an infinite amount of times), in accordance with the notion of a fair coin (which we encoded in our simulation by having `np.random.choice` choose between two possibilities with equal probability), with the fluctuations seeming to become much smaller as the number of trials increases.

Multiple Replications of the Coin Flip Experiment

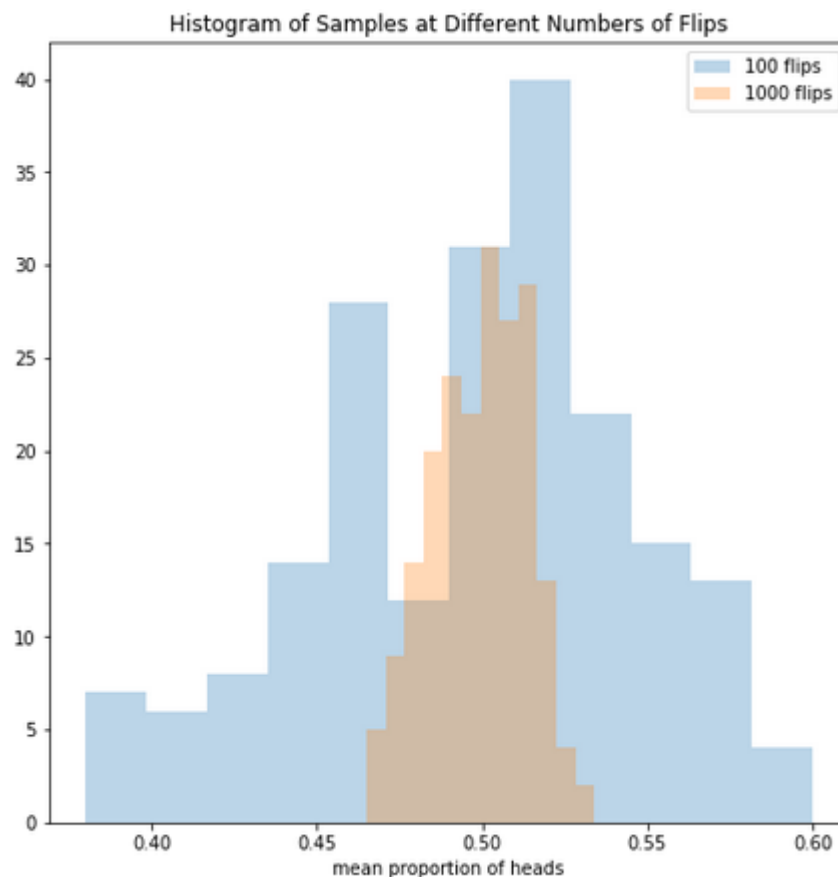
The coin flip experiment that we did above gave us some insight, but we don't have a good notion of how robust our results are under repetition as we've only run one experiment for each number of coin flips. Lets redo the coin flip experiment, but let's incorporate multiple repetitions of each

number of coin flips. For each choice of the number of flips, n , in an experiment, we'll do M replications of the coin tossing experiment.

Question 3. So Many Replications

3.1. Write a function `make_throws` which takes as arguments the `n_replications` (M) and the `n_flips` (n), and returns a list (of size M) of proportions, with each proportion calculated by taking the ratio of heads to total number of coin flips in each replication of n coin tosses. `n_flips` should be a python parameter whose value should default to 20 if unspecified when `make_throws` is called.

3.2. Create the variables `proportions_at_n_flips_100` and `proportions_at_n_flips_1000`. Store in these variables the result of `make_throws` for `n_flips` equal to 100 and 1000 respectively while keeping `n_replications` at 200. Create a plot with the histograms of `proportions_at_n_flips_100` and `proportions_at_n_flips_1000`. Make sure to title your plot, label the x-axis and provide a legend. (See below for an example of what the plot may look like)



3.3. Calculate the mean and variance of the results in each of the variables `proportions_at_n_flips_100` and `proportions_at_n_flips_1000` generated in 3.2.

3.4. Based upon the plots what would be your guess of what type of distribution is represented by histograms in 3.2? Explain the factors that influenced your choice.

- A. Gamma Distribution
- B. Beta Distribution
- C. Gaussian

3.5. Let's just assume for arguments sake that the answer to 3.4 is **C. Gaussian**. Plot a **normed histogram** of your results `proportions_at_n_flips_1000` overlayed with your selection for the appropriate gaussian distribution to represent the experiment of flipping a coin 1000 times. (**Hint: What parameters should you use for your Gaussian?**)

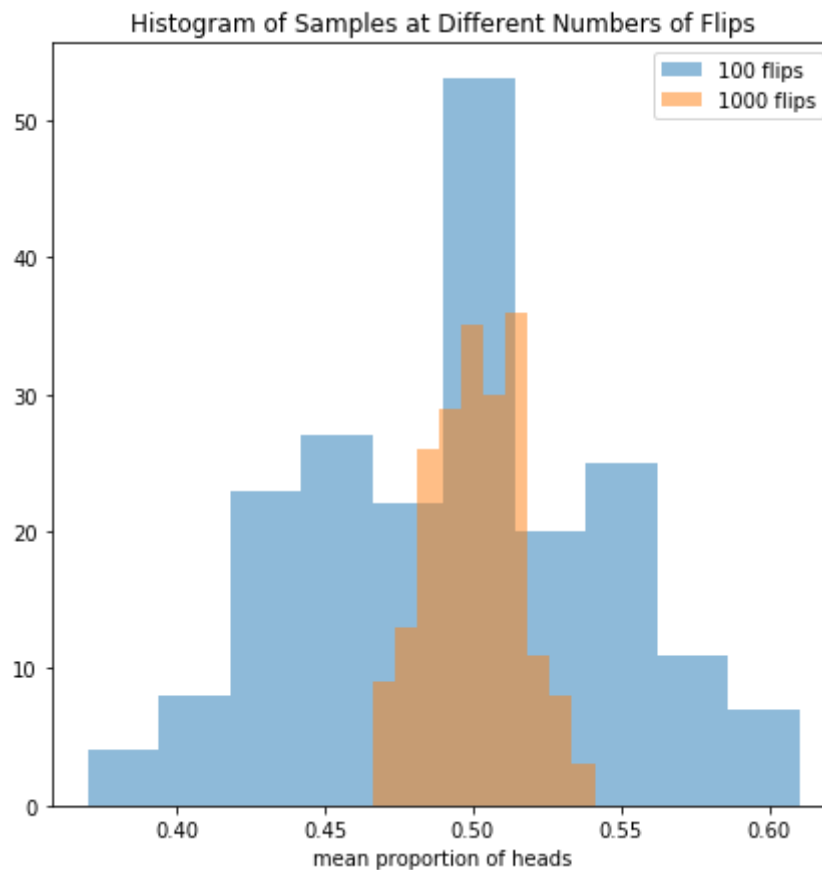
3.1

```
In [9]: # your code here
def make_throws(n_replications, n_flips = 20):
    return [find_prop_heads(throw_a_coin(n_flips), n_flips) for _ in range(n_replications)]
```

3.2

```
In [10]: # your code here
proportions_at_n_flips_100 = make_throws(200, 100)
proportions_at_n_flips_1000 = make_throws(200, 1000)
```

```
In [11]: # code for your plot here
plt.figure(figsize=(7, 7))
plt.hist(proportions_at_n_flips_100, alpha = 0.5, label = '100 flips')
plt.hist(proportions_at_n_flips_1000, alpha = 0.5, label = '1000 flips')
plt.xlabel('mean proportion of heads')
plt.title('Histogram of Samples at Different Numbers of Flips')
plt.legend(loc='upper right')
plt.show()
```



3.3

```
In [12]: # your code here
print('The mean at 100 flips is %f' % np.mean(proportions_at_n_flips_100))
print('The variance at 100 flips is %f' % np.var(proportions_at_n_flips_100))
print('The mean at 1000 flips is %f' % np.mean(proportions_at_n_flips_1000))
print('The variance at 1000 flips is %f' % np.var(proportions_at_n_flips_1000))
```

```
The mean at 100 flips is 0.493400
The variance at 100 flips is 0.002553
The mean at 1000 flips is 0.500660
The variance at 1000 flips is 0.000237
```

3.4

We can model the proportion of heads as the average of i.i.d. Bernoulli random variables with probability 0.5. By the central limit theorem, this average is approximately normally distributed as the

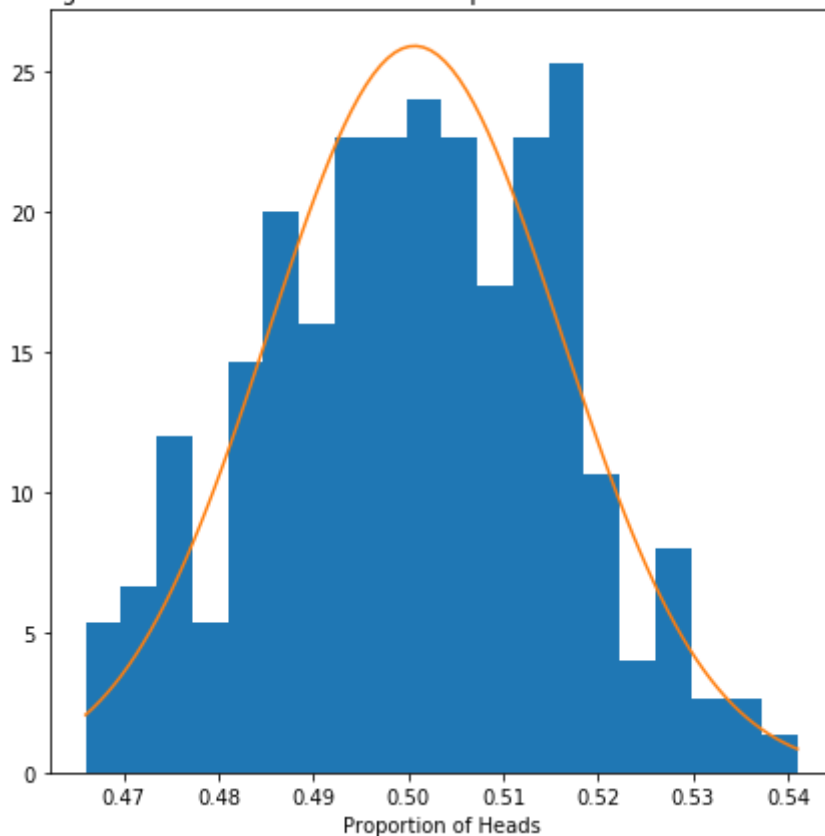
number of flips increases. Hence, the distribution of the proportion of heads is best approximated by a Gaussian distribution.

3.5

```
In [13]: # your code here
from scipy.stats import norm
mean = np.mean(proportions_at_n_flips_1000)
standard_dev = np.std(proportions_at_n_flips_1000)
x = np.linspace(np.min(proportions_at_n_flips_1000), np.max(proportions_at_n_flips_1000))
normal_pdf = norm.pdf(x, mean, standard_dev)
```

```
In [14]: plt.figure(figsize=(7, 7))
plt.hist(proportions_at_n_flips_1000, bins = 20, density = True)
plt.plot(x, normal_pdf)
plt.xlabel('Proportion of Heads')
plt.title('Fitting a Normal Distribution to the Proportion of Heads within 1000 Flips')
plt.show()
```

Fitting a Normal Distribution to the Proportion of Heads within 1000 Flips



Working With Distributions in Numpy/Scipy

Earlier in this problem set we've been introduced to the Bernoulli "aka coin-flip" distribution and worked with it indirectly by using `np.random.choice` to make a random selection between two elements 'H' and 'T'. Let's see if we can create comparable results by taking advantage of the machinery for working with other probability distributions in python using numpy and scipy.

Question 4: My Normal Binomial

Let's use our coin-flipping machinery to do some experimentation with the binomial distribution. The binomial distribution, often represented by $k \sim \text{Binomial}(n, p)$ is often described the number of successes in n Bernoulli trials with each trial having a probability of success p . In other words, if you flip a coin n times, and each coin-flip has a probability p of landing heads, then the number of heads you observe is a sample from a binomial distribution.

4.1. Sample the binomial distribution with $p = 0.5$ using coin flips by writing a function `sample_binomial1` which takes in integer parameters `n` and `size`. The output of `sample_binomial1` should be a list of length `size` observations with each observation being the outcome of flipping a coin `n` times and counting the number of heads. By default `size` should be 1. Your code should take advantage of the `throw_a_coin` function we defined above.

4.2. Sample the binomial distribution directly using `scipy.stats.binom.rvs` by writing another function `sample_binomial2` that takes in integer parameters `n` and `size` as well as a float `p` parameter `p` where $p \in [0..1]$. The output of `sample_binomial2` should be a list of length `size` observations with each observation a sample of $\text{Binomial}(n, p)$ (taking advantage of `scipy.stats.binom`). By default `size` should be 1 and `p` should be 0.5.

4.3. Run `sample_binomial1` with 25 and 200 as values of the `n` and `size` parameters respectively and store the result in `binomial_trials1`. Run `sample_binomial2` with 25, 200 and 0.5 as values of the `n`, `size` and `p` parameters respectively and store the results in `binomial_trials2`. Plot normed histograms of `binomial_trials1` and `binomial_trials2`. On both histograms, overlay a plot of the pdf of $\text{Binomial}(n = 25, p = 0.5)$

4.4. How do the plots in 4.3 compare?

4.5. Find the mean and variance of `binomial_trials1`. How do they compare to the true mean and variance of a $\text{Binomial}(n = 25, p = 0.5)$ distribution?

Answers

4.1

```
In [15]: def sample_binomial1(n, size = 1):
         return [np.sum(throw_a_coin(n) == 'H') for _ in range(size)]
```

4.2

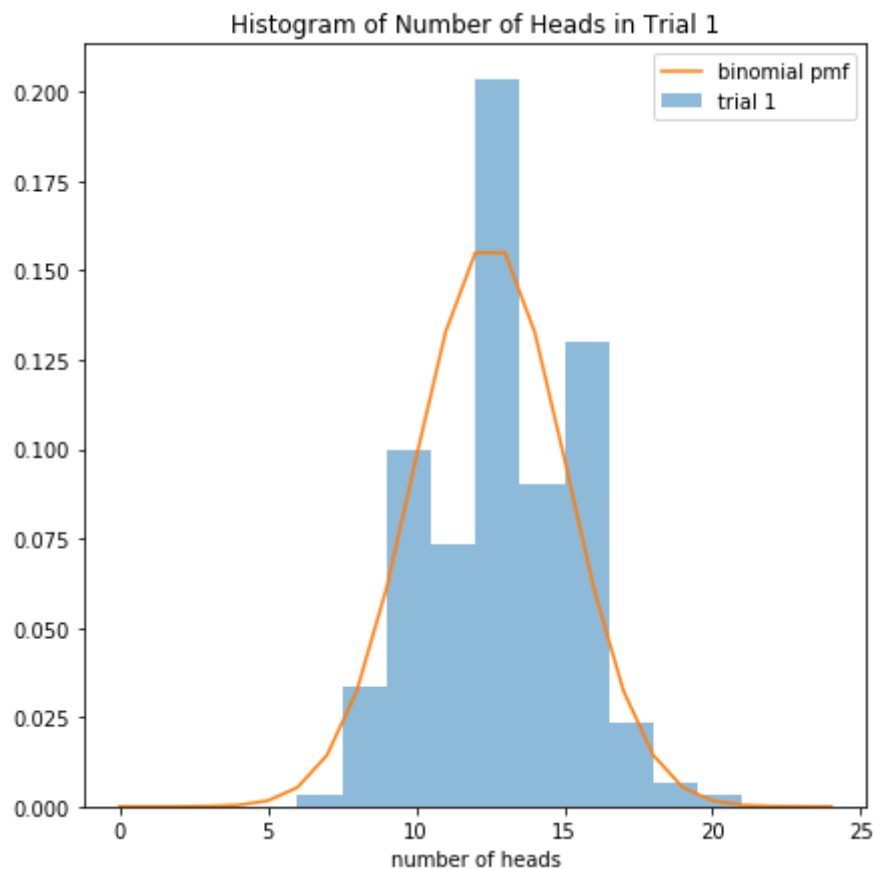
```
In [16]: # your code
         def sample_binomial2(n, p = 0.5, size = 1):
             return list(scipy.stats.binom.rvs(n, p, size = size))
```

4.3

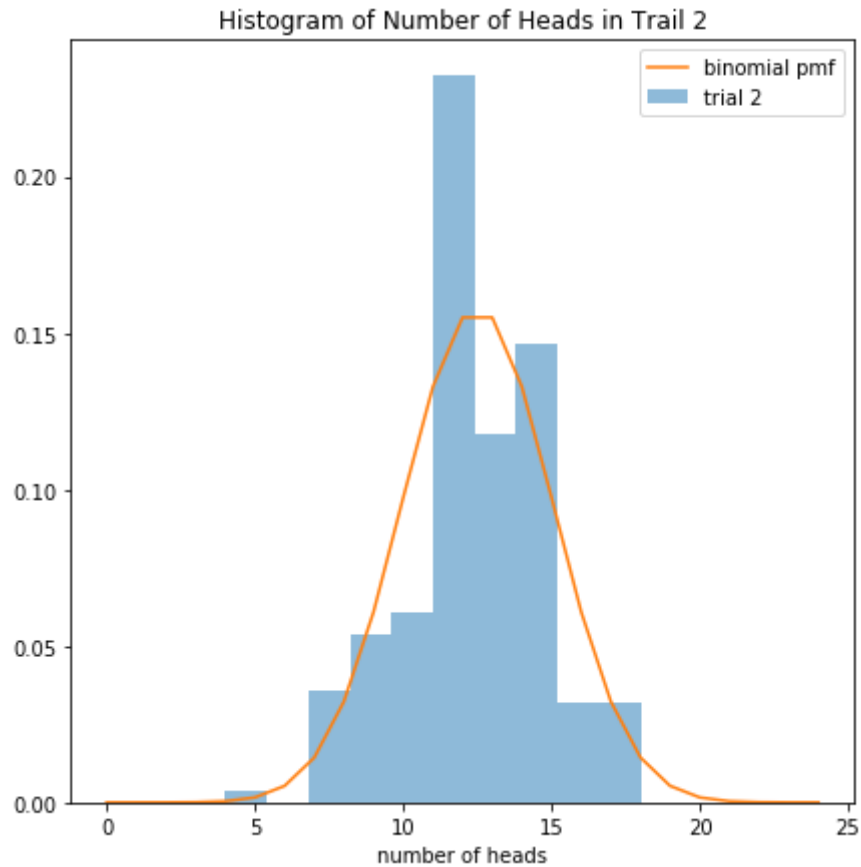

```
In [17]: # your code here
# your code here
binomial_trials1 = sample_binomial1(25, 200)
binomial_trials2 = sample_binomial2(25, 0.5, 200)
x = np.arange(0,25,1)
pmf = scipy.stats.binom.pmf(x, 25, 0.5)
pmf
```

```
Out[17]: array([2.98023224e-08, 7.45058060e-07, 8.94069672e-06, 6.85453415e-05,
 3.76999378e-04, 1.58339739e-03, 5.27799129e-03, 1.43259764e-02,
 3.22334468e-02, 6.08853996e-02, 9.74166393e-02, 1.32840872e-01,
 1.54981017e-01, 1.54981017e-01, 1.32840872e-01, 9.74166393e-02,
 6.08853996e-02, 3.22334468e-02, 1.43259764e-02, 5.27799129e-03,
 1.58339739e-03, 3.76999378e-04, 6.85453415e-05, 8.94069672e-06,
 7.45058060e-07])
```

```
In [18]: plt.figure(figsize=(7, 7))
plt.hist(binomial_trials1, alpha = 0.5, density = True, label = 'trial 1')
plt.plot(x, pmf, label = 'binomial pmf')
plt.xlabel('number of heads')
plt.title('Histogram of Number of Heads in Trial 1')
plt.legend(loc='upper right')
plt.show()
```



```
In [19]: plt.figure(figsize=(7, 7))
plt.hist(binomial_trials2, alpha = 0.5, density = True, label = 'trial 2')
plt.plot(x, pmf, label = 'binomial pmf')
plt.xlabel('number of heads')
plt.title('Histogram of Number of Heads in Trail 2')
plt.legend(loc='upper right')
plt.show()
```



4.4

Your explanation here

The binomial pmf fits the two trials well, although there is more variation in the 10-15 heads range.

4.5

```
In [20]: # your code here
print("The mean of the 1st trial is %f" % np.mean(binomial_trials1))
print("The variance of the 1st trial is %f" % np.var(binomial_trials1))
print("The mean of the binomial distribution with n = 25 and p = 0.5 is %f"
      % sp.stats.binom.stats(25, 0.5, moments = 'm'))
print("The variance of the binomial distribution with n = 25 and p = 0.5 is
      % sp.stats.binom.stats(25, 0.5, moments = 'v'))
```

```
The mean of the 1st trial is 12.630000
The variance of the 1st trial is 6.263100
The mean of the binomial distribution with n = 25 and p = 0.5 is 12.50000
0
The variance of the binomial distribution with n = 25 and p = 0.5 is 6.25
0000
```

**** Your explanation here ****

The mean and variance of the simulated data is close to the mean and variance of the binomial distribution with the appropriate parameters.

Testing Your Python Code

In the following section we're going to do a brief introduction to unit testing. We do so not only because unit testing has become an increasingly important part of the methodology of good software practices, but also because we plan on using unit tests as part of our own CS109 grading practices as a way of increasing rigor and repeatability decreasing complexity and manual workload in our evaluations of your code. We'll provide an example unit test at the end of this section.

Introduction to unit testing

```
In [21]: import ipytest
```

Unit testing is one of the most important software testing methodologies. Wikipedia describes unit testing as "a software testing method by which individual units of source code, sets of one or more computer program modules together with associated control data, usage procedures, and operating procedures, are tested to determine whether they are fit for use."

There are many different python libraries that support software testing in general and unit testing in particular. PyTest is one of the most widely used and well-liked libraries for this purpose. We've chosen to adopt PyTest (and ipytest which allows pytest to be used in ipython notebooks) for our testing needs and we'll do a very brief introduction to Pytest here so that you can become familiar with it too.

If you recall the function that we provided you above `throw_a_coin`, which we'll reproduce here for convenience, it took a number and returned that many "coin tosses". We'll start by seeing what happens when we give it different sizes of N . If we give $N = 0$, we should get an empty array of "experiments".

```
In [22]: def throw_a_coin(N):
         return np.random.choice(['H','T'], size=N)
```

```
In [23]: throw_a_coin(0)
```

```
Out[23]: array([], dtype='|S1')
```

Great! If we give it positive values of N we should get that number of 'H's and 'T's.

```
In [24]: throw_a_coin(5)
```

```
Out[24]: array(['T', 'T', 'H', 'H', 'H'], dtype='|S1')
```

```
In [25]: throw_a_coin(8)
```

```
Out[25]: array(['T', 'T', 'H', 'T', 'T', 'H', 'H', 'T'], dtype='|S1')
```

Exactly what we expected!

What happens if the input isn't a positive integer though?

```
In [26]: throw_a_coin(4.5)
```

```
-----
--
TypeError                                Traceback (most recent call last)
<ipython-input-26-7a98054470df> in <module>()
----> 1 throw_a_coin(4.5)

<ipython-input-22-9b62022d816e> in throw_a_coin(N)
      1 def throw_a_coin(N):
----> 2     return np.random.choice(['H','T'], size=N)

mtrand.pyx in mtrand.RandomState.choice()

mtrand.pyx in mtrand.RandomState.randint()

mtrand.pyx in mtrand.RandomState.randint()

randint_helpers.pxi in mtrand._rand_int64()

TypeError: 'float' object cannot be interpreted as an index
```

or

```
In [27]: throw_a_coin(-4)
```

```
-----
--
ValueError                                Traceback (most recent call las
t)
<ipython-input-27-8560c28a4e91> in <module>()
----> 1 throw_a_coin(-4)

<ipython-input-22-9b62022d816e> in throw_a_coin(N)
      1 def throw_a_coin(N):
----> 2     return np.random.choice(['H', 'T'], size=N)

mtrand.pyx in mtrand.RandomState.choice()

mtrand.pyx in mtrand.RandomState.randint()

mtrand.pyx in mtrand.RandomState.randint()

randint_helpers.pxi in mtrand._rand_int64()

ValueError: negative dimensions are not allowed
```

It looks like for both real numbers and negative numbers, we get two kinds of errors a `TypeError` and a `ValueError`. We just engaged in one of the most rudimentary forms of testing, trial and error. We can use `pytest` to automate this process by writing some functions that will automatically (and potentially repeatedly) test individual units of our code methodology. These are called **unit tests**.

Before we write our tests, let's consider what we would think of as the appropriate behavior for `throw_a_coin` under the conditions we considered above. If `throw_a_coin` receives positive integer input, we want it to behave exactly as it currently does -- returning an output consisting of a list of characters 'H' or 'T' with the length of the list equal to the positive integer input. For a positive floating point input, we want `throw_a_coin_properly` to treat the input as if it were rounded down to the nearest integer thus returning a list of 'H' or 'T' integers whose length is the same as the input rounded down to the next highest integer. For a any negative number input or an input of 0, we want `throw_a_coin_properly` to return an empty list.

We create `pytest` tests by writing functions that start or end with "test". We'll use the **convention** that our tests will start with "test".

We begin the code cell with `ipytest`'s `clean_tests` function as a way to clear out the results of previous tests starting with "test_throw_a_coin" (the * is the standard wild card charater here).

```
In [28]: ## the * after test_throw_a_coin tells this code cell to clean out the results
## of all tests starting with test_throw_a_coin
ipytest.clean_tests("test_throw_a_coin*")

## run throw_a_coin with a variety of positive integer inputs (all numbers
## verify that the length of the output list (e.g ['H', 'H', 'T', 'H', 'T'])
def test_throw_a_coin_length_positive():
    for n in range(1,20):
        assert len(throw_a_coin(n)) == n

## verify that throw_a_coin produces an empty list (i.e. a list of length 0
## of 0
def test_throw_a_coin_length_zero():
    ## should be the empty array
    assert len(throw_a_coin(0)) == 0

## verify that given a positive floating point input (i.e. 4.34344298547201
## coin flips of length equal to highest integer less than the input
def test_throw_a_coin_float():
    for n in np.random.exponential(7, size=5):
        assert len(throw_a_coin(n)) == np.floor(n)

## verify that given any negative input (e.g. -323.4), throw_a_coin produces
def test_throw_a_coin_negative():
    for n in range(-7, 0):
        assert len(throw_a_coin(n)) == 0

ipytest.run_tests()
```

```
unittest.case.FunctionTestCase (test_throw_a_coin_float) ... ERROR
unittest.case.FunctionTestCase (test_throw_a_coin_length_positive) ... ok
unittest.case.FunctionTestCase (test_throw_a_coin_length_zero) ... ok
unittest.case.FunctionTestCase (test_throw_a_coin_negative) ... ERROR
```

```
=====
ERROR: unittest.case.FunctionTestCase (test_throw_a_coin_float)
-----
```

```
Traceback (most recent call last):
```

```
File "<ipython-input-28-78a86d656b91>", line 22, in test_throw_a_coin_float
    assert len(throw_a_coin(n)) == np.floor(n)
```

```
File "<ipython-input-22-9b62022d816e>", line 2, in throw_a_coin
```

```
    return np.random.choice(['H','T'], size=N)
```

```
File "mtrand.pyx", line 1163, in mtrand.RandomState.choice
```

```
File "mtrand.pyx", line 995, in mtrand.RandomState.randint
```

```
File "mtrand.pyx", line 996, in mtrand.RandomState.randint
```

```
File "randint_helpers.pxi", line 253, in mtrand._rand_int64
```

```
TypeError: 'numpy.float64' object cannot be interpreted as an index
```

```
=====
ERROR: unittest.case.FunctionTestCase (test_throw_a_coin_negative)
-----
```

```
Traceback (most recent call last):
```

```
File "<ipython-input-28-78a86d656b91>", line 28, in test_throw_a_coin_n
```

```

negative
    assert len(throw_a_coin(n)) == 0
File "<ipython-input-22-9b62022d816e>", line 2, in throw_a_coin
    return np.random.choice(['H','T'], size=N)
File "mtrand.pyx", line 1163, in mtrand.RandomState.choice
File "mtrand.pyx", line 995, in mtrand.RandomState.randint
File "mtrand.pyx", line 996, in mtrand.RandomState.randint
File "randint_helpers.pxi", line 253, in mtrand._rand_int64
ValueError: negative dimensions are not allowed

```

```

-----
Ran 4 tests in 0.009s

```

```

FAILED (errors=2)

```

As you see, we were able to use pytest (and ipytest which allows us to run pytest tests in our ipython notebooks) to automate the tests that we constructed manually before and get the same errors and successes. Now time to fix our code and write our own test!

Question 5: You Better Test Yourself before You Wreck Yourself!

Now it's time to fix `throw_a_coin` so that it passes the tests we've written above as well as add our own test to the mix!

5.1. Write a new function called `throw_a_coin_properly` that will pass the tests that we saw above. For your convenience we'll provide a new jupyter notebook cell with the tests rewritten for the new function. All the tests should pass. For a positive floating point input, we want `throw_a_coin_properly` to treat the input as if it were rounded down to the nearest integer. For any negative number input, we want `throw_a_coin_properly` to treat the input as if it were 0.

5.2. Write a new test for `throw_a_coin_properly` that verifies that all the elements of the resultant arrays are 'H' or 'T'.

Answers

5.1

```

In [29]: # your code here
def throw_a_coin_properly(n_trials):
    if isinstance(n_trials, float):
        n_trials = np.floor(n_trials).astype(int)
    if n_trials < 0:
        n_trials = 0
    return np.random.choice(['H','T'], size=n_trials)

```

```

In [30]: for n in np.random.exponential(7, size=5):
    assert len(throw_a_coin_properly(n)) == np.floor(n)

```

```
In [31]: ipytest.clean_tests("test_throw_a_coin*")

def test_throw_a_coin_properly_length_positive():
    for n in range(1,20):
        assert len(throw_a_coin_properly(n)) == n

def test_throw_a_coin_properly_length_zero():
    ## should be the empty array
    assert len(throw_a_coin_properly(0)) == 0

def test_throw_a_coin_properly_float():

    for n in np.random.exponential(7, size=5):
        assert len(throw_a_coin_properly(n)) == np.floor(n)

def test_throw_a_coin_properly_negative():

    for n in range(-7, 0):
        assert len(throw_a_coin_properly(n)) == 0

ipytest.run_tests()

unittest.case.FunctionTestCase (test_throw_a_coin_properly_float) ... ok
unittest.case.FunctionTestCase (test_throw_a_coin_properly_length_positiv
e) ... ok
unittest.case.FunctionTestCase (test_throw_a_coin_properly_length_zero)
... ok
unittest.case.FunctionTestCase (test_throw_a_coin_properly_negative) ...
ok

-----
Ran 4 tests in 0.007s

OK
```

5.2


```
In [32]: ipytest.clean_tests("test_throw_a_coin*")

## write a test that verifies you don't have any other elements except H's
def test_throw_a_coin_properly_verify_H_T():
    n = 20
    for toss in throw_a_coin_properly(n):
        assert toss == "H" or toss == "T"
    # your code here

ipytest.run_tests()
```

```
unittest.case.FunctionTestCase (test_throw_a_coin_properly_verify_H_T)
... ok
```

```
-----
Ran 1 test in 0.002s
```

OK

```
In [33]: from IPython.core.display import HTML
def css_styling():
    styles = open("cs109.css", "r").read()
    return HTML(styles)
css_styling()
```

Out[33]: