

# 3038741162\_Workshop05

March 17, 2023

## 1 Workshop 5: PDF sampling and Statistics

Submit this notebook to bCourses to receive a grade for this Workshop.

Please complete workshop activities in code cells in this iPython notebook. The activities titled **Practice** are purely for you to explore Python, and no particular output is expected. Some of them have some code written, and you should try to modify it in different ways to understand how it works. Although no particular output is expected at submission time, it is *highly* recommended that you read and work through the practice activities before or alongside the exercises. However, the activities titled **Exercise** have specific tasks and specific outputs expected. Include comments in your code when necessary.

The workshop should be submitted on bCourses under the Assignments tab (both the .ipynb and .pdf files).

### 1.1 Preview: generating random numbers

We will discuss simulations in greater detail later in the semester. The first step in simulating nature – which, despite Einstein’s objections, is playing dice after all – is to learn how to generate some numbers that appear random. Of course, computers cannot generate true random numbers – they have to follow an algorithm. But the algorithm may be based on something that is difficult to predict (e.g. the time of day you are executing this code) and therefore *look* random to a human. Sequences of such numbers are called *pseudo-random*.

The random variables you generate will be distributed according to some *Probability Density Function* (PDF). The most common PDF is *flat*:  $f(x) = \frac{1}{b-a}$  for  $x \in [a..b)$ . Here is how to get a random number uniformly distributed between  $a = 0$  and  $b = 1$  in Python:

```
[2]: # standard preamble
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

```
[3]: # generate one random number between [0,1)
x = np.random.rand()
print ('x=', x)

# generate an array of 10 random numbers between [0,1)
array = np.random.rand(10)
print (array)
```

```
x= 0.6996036343917029
[0.22126755 0.90226071 0.33845573 0.27653026 0.86423314 0.67955987
 0.55847681 0.82137739 0.4980035 0.28297366]
```

You can generate a set of randomly-distributed integer values instead:

```
[5]: a = np.random.randint(0,1000,10)
      print(a)
```

```
[763 179 95 526 47 386 127 610 398 659]
```

## 2 1d distributions

### 2.1 Moments of the distribution

Python's SciPy library contains a set of standard statistical functions. See a few examples below:

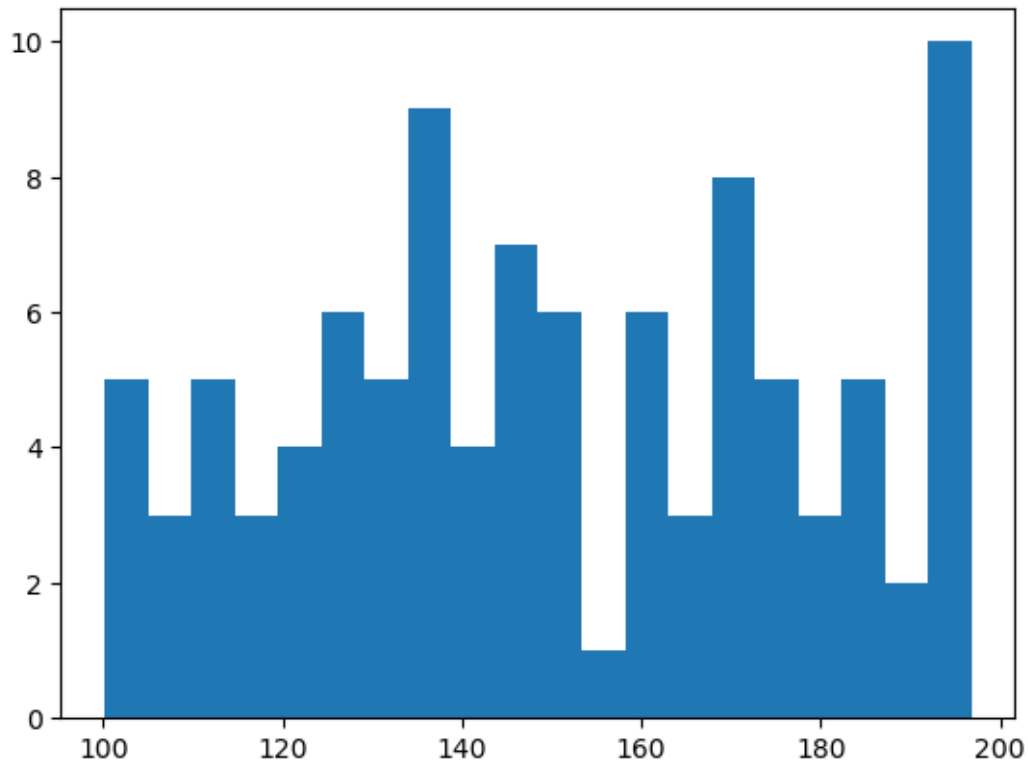
```
[9]: # create a set of data and compute mean and variance
      # This creates an array of 100 elements, uniformly-distributed between 100 and 200
      # Try changing the size parameter!
      x = np.random.uniform(low=100,high=200,size=100)

      print(x[0:10])
      # make a histogram
      n, bins, patches = plt.hist(x, 20)
      #20: bins(      =      )

      # various measures of "average value":
      print('Mean = {0:5.0f}'.format(np.mean(x)))
      print('Median = {0:5.0f}'.format(np.median(x)))

      # measure of the spread
      print('Standard deviation = {0:5.1f}'.format(np.std(x)))
```

```
[109.00420264 182.99138173 150.65889297 171.34541982 169.34764919
 192.99253398 110.42760158 125.67085025 185.40923449 139.60652038]
Mean =    150
Median =    147
Standard deviation =   27.8
```



### 2.1.1 Exercise 1

We just introduced some new functions: `np.random.rand()`, `np.random.uniform()`, `plt.hist()`, `np.mean()`, and `np.median()`. So let's put them to work. You may also find `np.cos()`, `np.sin()`, and `np.std()` useful.

1. Generate 100 random numbers, uniformly distributed between  $[-\pi, \pi)$
2. Plot them in a histogram.
3. Compute mean and standard deviation (RMS)
4. Plot a histogram of  $\sin(x)$  and  $\cos(x)$ , where  $x$  is a uniformly distributed random number between  $[-\pi, \pi)$ . Do you understand this distribution ?

```
[44]: # Your code for Exercise 1
#1. Generate 100 random numbers, uniformly distributed between [-, )
x = np.random.uniform(low = -np.pi, high = np.pi, size = 100)

fig, ax = plt.subplots(2, 2, figsize=(9, 7))
ax[0,0].hist(x, 30)
ax[0,0].set_title("x")

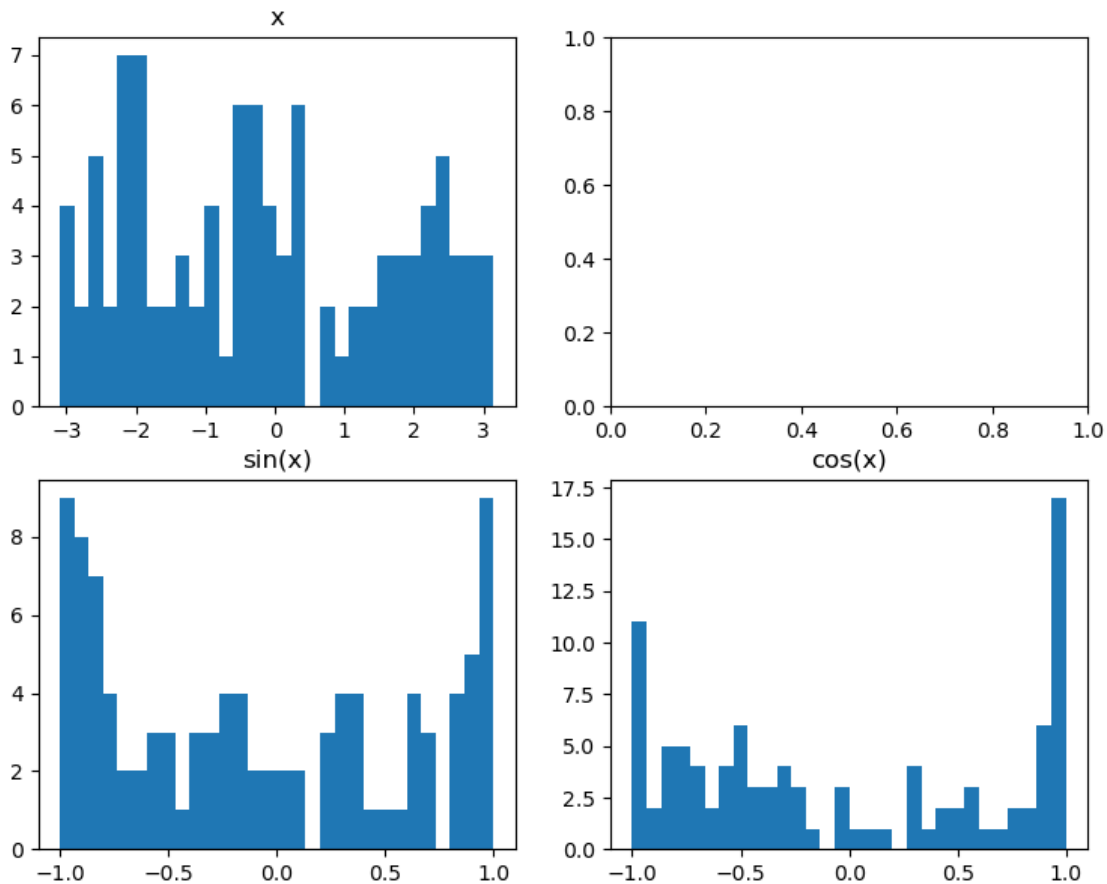
ax[1,0].hist(np.sin(x), 30)
ax[1,0].set_title("sin(x)")
```

```
ax[1,1].hist(np.cos(x), 30)
ax[1,1].set_title("cos(x)")

print('Mean = {0:5.5f}'.format(np.mean(x)))
print('Standard deviation = {0:5.5f}'.format(np.std(x)))
```

Mean = -0.13863

Standard deviation = 1.83502



## 2.2 Gaussian/Normal distribution

You can also generate Gaussian-distributed numbers. Remember that a Gaussian (or Normal) distribution is a probability distribution given by

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

where  $\mu$  is the average of the distribution and  $\sigma$  is the standard deviation. The **standard** normal distribution is a special case with  $\mu = 0$  and  $\sigma = 1$ .

```
[45]: # generate a single random number, gaussian-distributed with mean=0 and sigma=1.
      ↪ This is also called
      # a standard normal distribution
      x = np.random.standard_normal()
      print (x)

      # generate an array of 10 such numbers
      a = np.random.standard_normal(size=10)
      print (a)
```

```
-0.1836190546275852
[ 0.21962384 -2.58416995  1.31989456  0.07069773 -0.51734238 -0.13619946
 -0.20579235  0.25421833  0.76213769  0.42766107]
```

### 2.2.1 Exercise 2

We now introduced `np.random.standard_normal()`.

1. Generate  $N = 100$  random numbers, Gaussian-distributed with  $\mu = 0$  and  $\sigma = 1$ .
2. Plot them in a histogram.
3. Compute the mean, standard deviation (RMS), and standard error on the mean.

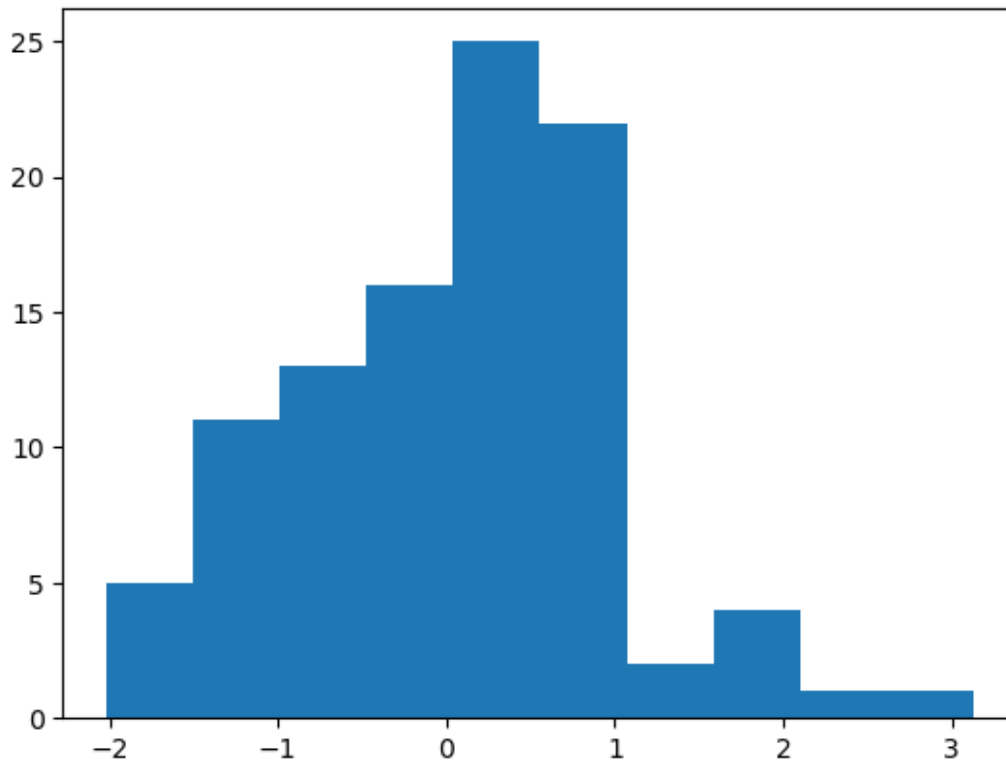
The standard error on the mean is defined as  $\sigma_\mu = \frac{\sigma}{\sqrt{N}}$ , where  $\sigma$  is the standard deviation.

```
[14]: # Your code for Exercise 2
      N = 100
      gauss = np.random.standard_normal(size = N)

      plt.hist(gauss, 10)

      print('Mean = {0:5.5f}'.format(np.mean(gauss)))
      std_deviation = np.std(gauss)
      print('Standard deviation = {0:5.5f}'.format(std_deviation))
      print('Standard error on the mean = {0:5.5f}'.format(std_deviation / np.
      ↪sqrt(N)))
```

```
Mean = 0.03988
Standard deviation = 0.94135
Standard error on the mean = 0.09413
```



4. Now find the means of  $M = 1000$  experiments of  $N = 100$  measurements each (you'll end up generating 100,000 random numbers total). Plot a histogram of the means. Is it consistent with your calculation of the error on the mean for  $N = 100$  ? About how many experiments yield a result within  $1\sigma_\mu$  of the true mean of 0 ? About how many are within  $2\sigma_\mu$  ?
5. Now repeat question 4 for  $N = 10, 50, 1000, 10000$ . Plot a graph of the RMS of the distribution of the means vs  $N$ . Is it consistent with your expectations ?

```
[75]: # 4
means = []
N = 100
M = 1000
for i in range(M):
    x = np.random.standard_normal(size = N)
    means.append(np.mean(x))

plt.hist(means, 100)
plt.show()

mean_of_means = np.mean(means)
std_error_of_mean = np.std(means, ddof=1) / np.sqrt(M)

print('Mean of means = {0:5.5f}'.format(mean_of_means))
```

```

print('Standard error on the mean = {0:5.5f}'.format(std_error_of_mean))

num_within_1se = np.sum(np.abs(means - mean_of_means) < std_error_of_mean)
num_within_2se = np.sum(np.abs(means - mean_of_means) < 2*std_error_of_mean)

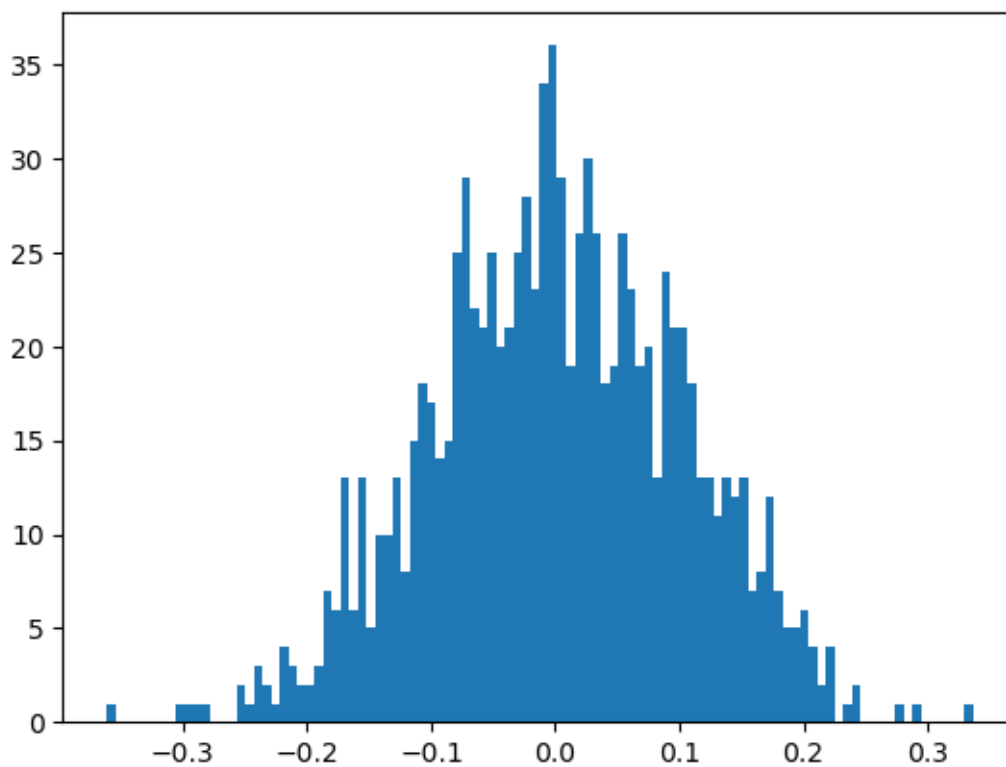
print('Number of experiments within 1 standard error of the mean = {0:d}'.
      ↪format(num_within_1se))
print('Number of experiments within 2 standard errors of the mean = {0:d}'.
      ↪format(num_within_2se))

# 5
Ns = [10, 50, 1000, 10000]
rms_list = []

for N in Ns:
    means = []
    for i in range(M):
        x = np.random.standard_normal(size=N)
        means.append(np.mean(x))
    rms = np.sqrt(np.mean(np.square(means)))
    rms_list.append(rms)

plt.plot(Ns, rms_list, '-o')
plt.xscale('log')
plt.xlabel('N')
plt.ylabel('RMS of means')
plt.show()

```



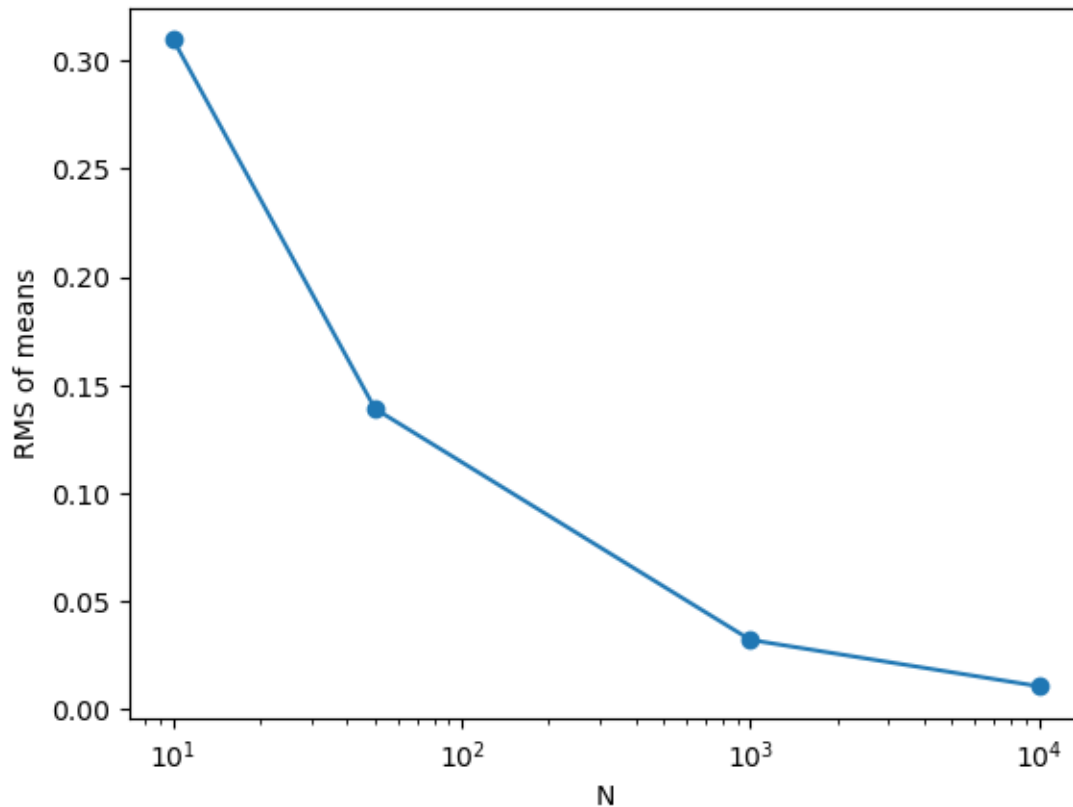
Mean of means = 0.00183

Standard error on the mean = 0.00322

Number of experiments within 1 standard error of the mean = 29

Number of experiments within 2 standard errors of the mean = 62





## 2.3 Exponential distribution

In this part we will repeat the above process, but now using lists of exponentially distributed random numbers. The probability of selecting a random number between  $x$  and  $x + dx$  is  $\propto e^{-x} dx$ . Exponential distributions often appear in lossy systems, e.g. if you plot an amplitude of a damped oscillator as a function of time. Or you may see it when you plot the number of decays of a radioactive isotope as a function of time.

```
[65]: # generate a single random number, exponentially-distributed with scale=1.
x = np.random.exponential()
print (x)

# generate an array of 10 such numbers
a = np.random.exponential(size=10)
print (a)
```

```
0.6194094348445207
[0.23777691 0.78297806 0.15130362 0.15045145 1.54514722 2.22356288
 0.77863909 0.9413876 0.61210442 0.06681621]
```

### 2.3.1 Exercise 3

We now introduced `np.random.exponential()`. This function can take up to two keywords, one of which is `size` as shown above. The other is `scale`. Use the documentation and experiment with this exercise to see what it does.

1. What do you expect to be the mean of the distribution? What do you expect to be the standard deviation?
2. Generate  $N = 100$  random numbers, exponentially-distributed with the keyword `scale` set to 1.
3. Plot them in a histogram.
4. Compute mean, standard deviation (RMS), and the error on the mean. Is this what you expected?
5. Now find the means, standard deviations, and errors on the means for each of the  $M = 1000$  experiments of  $N = 100$  measurements each. Plot a histogram of each quantity. Is the RMS of the distribution of the means consistent with your calculation of the error on the mean for  $N = 100$  ?
6. Now repeat question 5 for  $N = 10, 100, 1000, 10000$ . Plot a graph of the RMS of the distribution of the means vs  $N$ . Is it consistent with your expectations ? This is a demonstration of the *Central Limit Theorem*

```
[16]: # Your code for Exercise 3
N = 100
x = np.random.exponential(size = N, scale = 1)

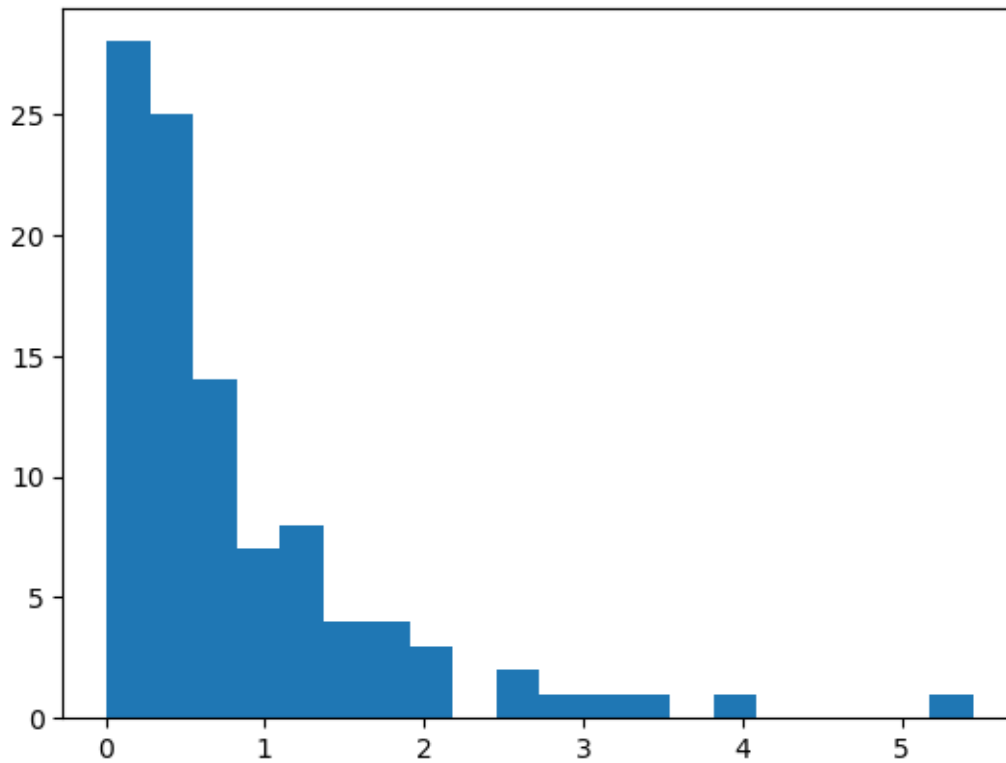
plt.hist(x, 20)

print('Mean = {0:5.5f}'.format(np.mean(x)))
std_deviation = np.std(x)
print('Standard deviation = {0:5.5f}'.format(std_deviation))
print('Standard error on the mean = {0:5.5f}'.format(std_deviation / np.
↪sqrt(N)))
```

Mean = 0.82584

Standard deviation = 0.92922

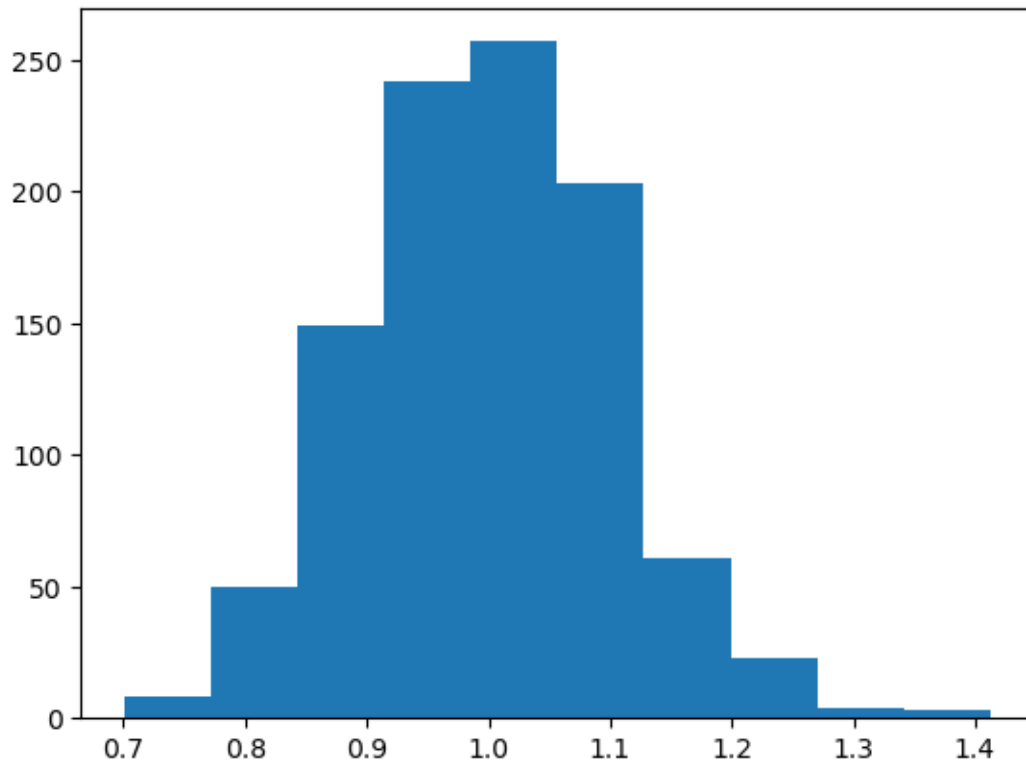
Standard error on the mean = 0.09292



```
[24]: #4
means = []
M = 1000

N = 100
for i in range(M):
    x = np.random.exponential(size = N, scale = 1)
    means.append(np.mean(x))

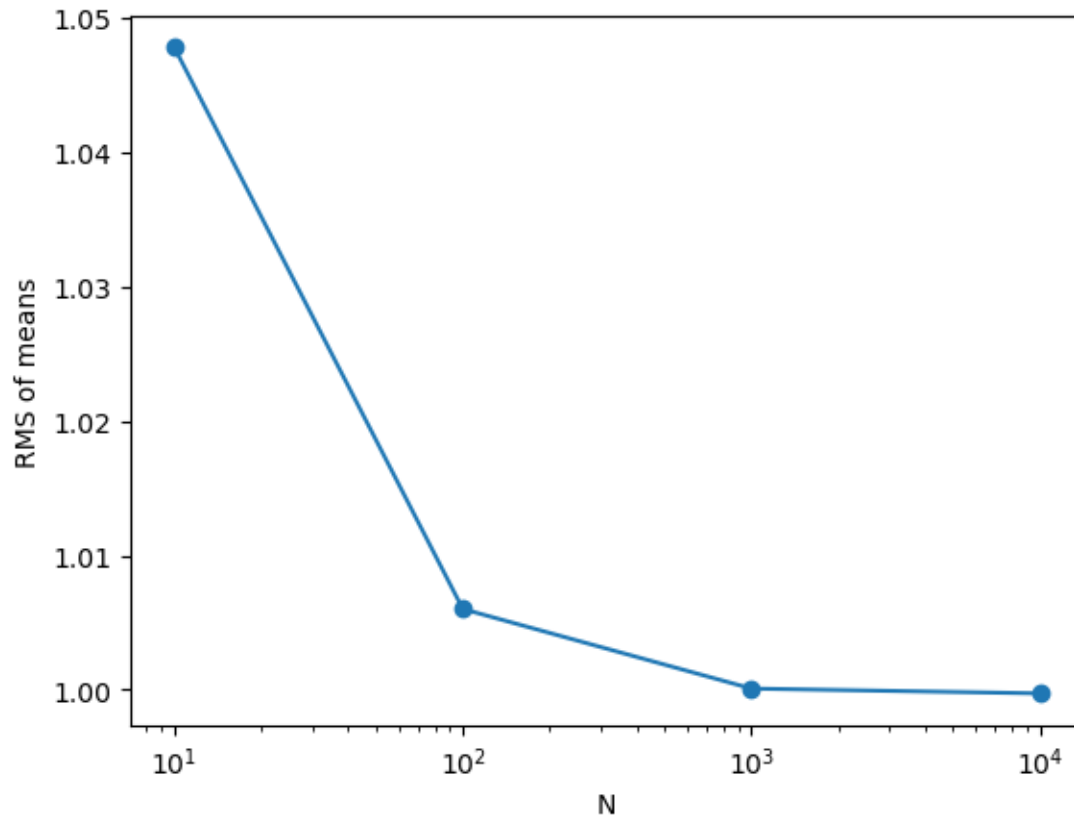
plt.hist(means)
plt.show()
```



```
[74]: #5
Ns = [10, 100, 1000, 10000]
rms_list = []

for N in Ns:
    means = []
    for i in range(M):
        x = np.random.exponential(size=N)
        means.append(np.mean(x))
    rms = np.sqrt(np.mean(np.square(means)))
    rms_list.append(rms)

plt.plot(Ns, rms_list, '-o')
plt.xscale('log')
plt.xlabel('N')
plt.ylabel('RMS of means')
plt.show()
```



## 2.4 Binomial distribution

The binomial distribution with parameters  $n$  and  $p$  is the *discrete* probability distribution of the number of successes in a sequence of  $n$  independent yes/no experiments, each of which yields success with probability  $p$ . A typical example is a distribution of the number of *heads* for  $n$  coin flips ( $p = 0.5$ )

```
[37]: # Simulates flipping 1 fair coin one time. Returns 0 for heads and 1 for tails
p = 0.5
print (np.random.binomial(1, p))

# Simulates flipping 5 biased coins three times
p = 0.7
print (np.random.binomial(5, p, size=3))
```

```
1
[3 4 3]
```

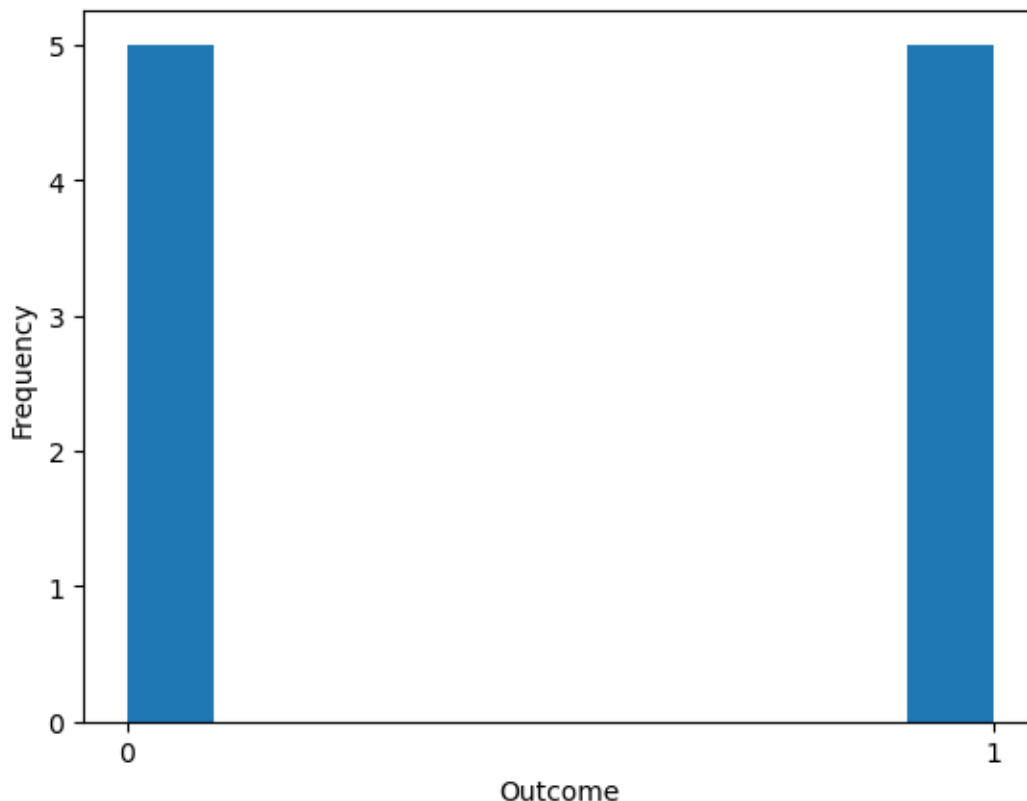
### 2.4.1 Exercise 4

We now introduced the function `np.random.binomial(n,p)` which requires two arguments, `n` the number of coins being flipped in a single trial and `p` the probability that a particular coin lands tails. As usual, `size` is another optional keyword argument. 1. Generate an array of outcomes for flipping 1 unbiased coin 10 times. 1. Plot the outcomes in a histogram (0=heads, 1=tails). 1. Compute mean, standard deviation (RMS), and the error on the mean. Is this what you expected?

```
[58]: # Your code for Exercise 4
N = 10
outcomes = np.random.binomial(1, 0.5, size=N)

plt.hist(outcomes)
#plt.hist(outcomes[0], bins=[-0.5, 0.5, 1.5], align='mid')
plt.xticks([0, 1])
plt.xlabel('Outcome')
plt.ylabel('Frequency')
plt.show()

print('Mean = {0:5.5f}'.format(np.mean(outcomes)))
std_deviation = np.std(outcomes)
print('Standard deviation = {0:5.5f}'.format(std_deviation))
print('Standard error on the mean = {0:5.5f}'.format(std_deviation / np.
    ↳sqrt(N)))
```



```
Mean = 0.50000
Standard deviation = 0.50000
Standard error on the mean = 0.15811
```

## 2.5 Poisson distribution

The Poisson distribution is a *discrete* probability distribution that expresses the probability of a given number of events  $n$  occurring in a fixed interval of time  $T$  if these events occur with a known average rate  $\nu/T$  and independently of the time since the last event. The *expectation value* of  $n$  is  $\nu$ . The variance of  $n$  is also  $\nu$ , so the standard deviation of  $n$  is  $\sigma(n) = \sqrt{\nu}$

```
[61]: nu = 10 # expected number of events
      n = np.random.poisson(nu) # generate a Poisson-distributed number.
      print (n)
```

13

### 2.5.1 Exercise 5

We introduced `np.random.poisson()`. As usual, you can use the keyword argument `size` to draw multiple samples. 1. Generate  $N = 100$  random numbers, Poisson-distributed with  $\nu = 10$ . 1. Plot them in a histogram. 1. Compute mean, standard deviation (RMS), and the error on the mean. Is this what you expected? 1. Now repeat question 3 for  $\nu = 1, 5, 100, 10000$ . Plot a graph of the RMS vs  $\nu$ . Is it consistent with your expectations ?

```
[73]: # Your code for Exercise 5
      N = 100
      nu = 10
      samples = np.random.poisson(nu, size=N)

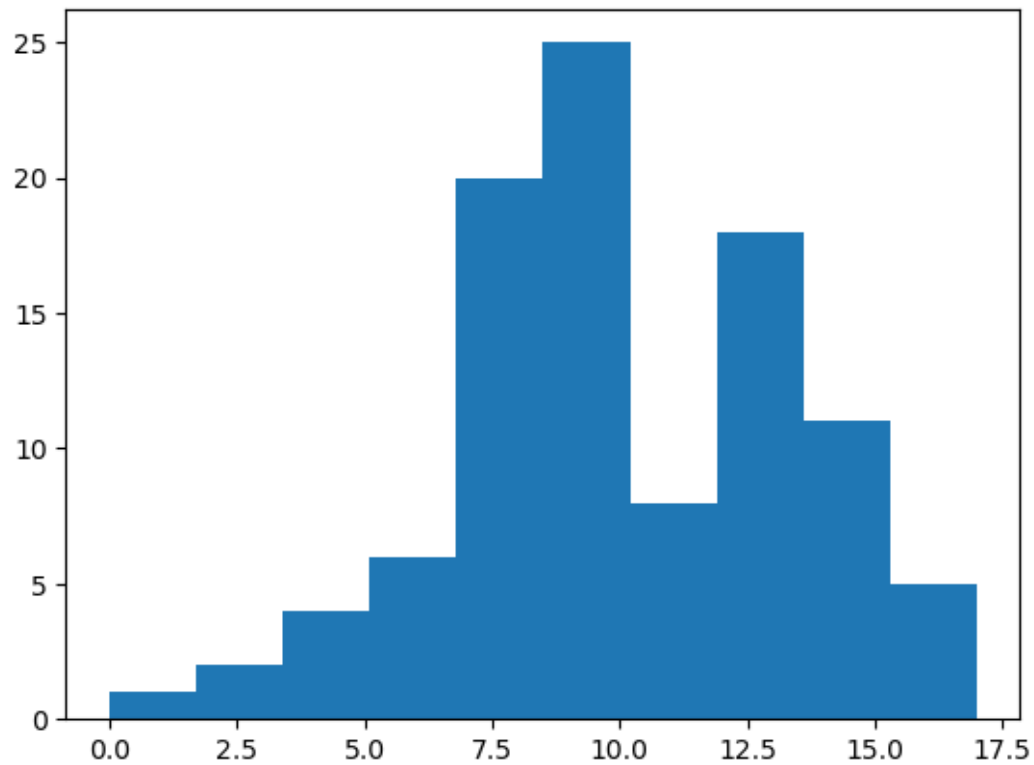
      plt.hist(samples, 10)
      plt.show()

      print('Mean = {0:5.5f}'.format(np.mean(samples)))
      std_deviation = np.std(samples)
      print('Standard deviation = {0:5.5f}'.format(std_deviation))
      print('Standard error on the mean = {0:5.5f}'.format(std_deviation / np.
        ↪sqrt(N)))

      Ns = [1, 5, 100, 10000]
      rms_list = []

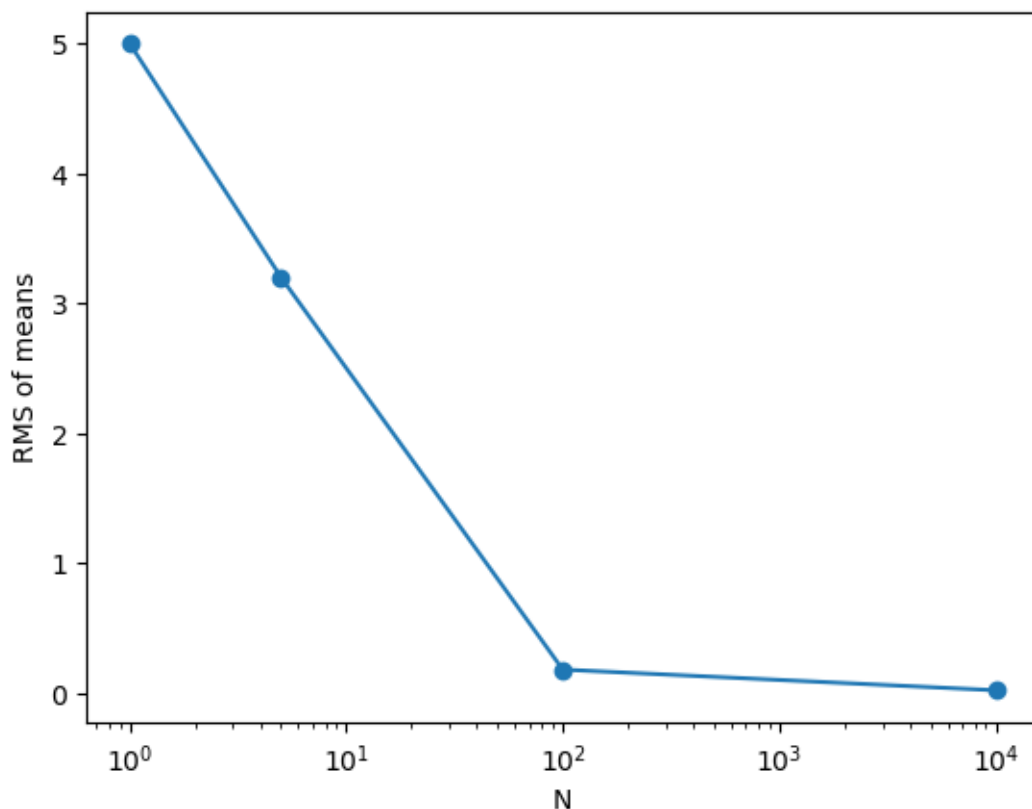
      for N in Ns:
          x = np.random.poisson(nu, size=N)
          rms = np.sqrt(np.mean(np.square(np.mean(x) - nu)))
          rms_list.append(rms)
```

```
plt.plot(Ns, rms_list, '-o')  
plt.xscale('log')  
plt.xlabel('N')  
plt.ylabel('RMS of means')  
plt.show()
```



Mean = 10.05000  
Standard deviation = 3.28139  
Standard error on the mean = 0.32814





## 2.6 Doing something “useful” with a distribution

[Random walks](#) show up when studying statistical mechanics (and many other fields). The simplest random walk is this:

Imagine a person stuck walking along a straight line. Each second, they randomly step either 1 meter forward or 1 meter backward.

With this in mind, you can start to ask many different questions. After one minute, how far do they end up from their starting point? How many times do they cross the starting point? (The exact answers require repeating this “experiment” many times and taking an average across all the trials.) How much do you have to pay someone to walk along this line for several hours?

There are lots of interesting ways to generalize this problem. You can extend the random walk to 2+ dimensions, make stepping in some directions more likely than others, draw the step sizes from some probability distribution, etc. If you’re curious, it’s fun to plot the paths of 2D random walks to visualize Brownian motion.

### 2.6.1 Exercise 6

Use `np.random.binomial(1, 0.5)` (or some other random number generator) to simulate a random walk along one dimension (the numbers from the binomial distribution signify either stepping

forward or backward). It would be helpful to write a function that takes  $N$  steps in the random walk, and then returns the distance from the starting point.

```
[81]: def random_walk(N):  
  
    '''This function will return the distance from the starting point  
        after a 1-dimensional random walk of  $N$  steps'''  
  
    # Use np.random.binomial(1,0.5) or another np.random function to "simulate"  
    ↪the random walk  
    steps = np.random.binomial(1, 0.5, size=N)  
    position = 0  
    for step in steps:  
        if step == 0:  
            position -= 1  
        else:  
            position += 1  
    return position
```

Now that you have a function that simulates a single random walk for a given  $N$ , write a function (or just some lines of code) that simulates  $M = 1000$  of these random walks and returns the mean (average) distance traveled for a given  $N$ .

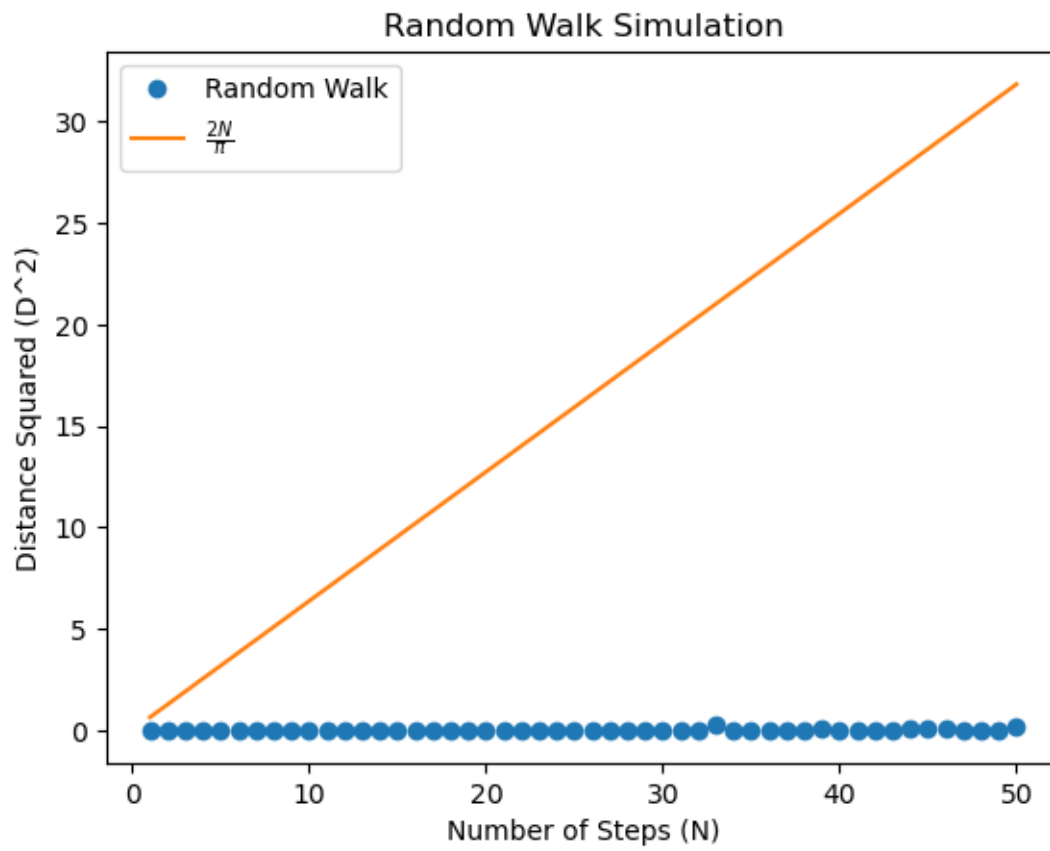
```
[82]: def average_distance(N):  
  
    '''This function simulates 1000 random walks of  $N$  steps  
        and then returns the average distance from the start.'''  
  
    # Use the random_walk(N) function 1000 times and return the average of the  
    ↪results  
    distances = []  
    for i in range(1000):  
        distances.append(random_walk(N))  
    return np.mean(distances)
```

It turns out that you can now use these random walk simulations to estimate the value of  $\pi$  (although in an extremely inefficient way). For values of  $N$  from 1 to 50, use your functions/code to find the mean distance  $D$  after  $N$  steps. Then make a plot of  $D^2$  vs  $N$ . If you've done it correctly, the plot should be a straight line with a slope of  $\frac{2}{\pi}$ .

Once we get to fitting in Python, you could find the slope and solve for  $\pi$ . For now, just draw the line  $\frac{2N}{\pi}$  over your simulated data.

```
[86]: N_vals = range(1, 51)  
D_vals = [average_distance(N)**2 for N in N_vals]  
  
plt.plot(N_vals, D_vals, 'o', label='Random Walk')  
plt.plot(N_vals, [(2*N)/np.pi for N in N_vals], label=r'$\frac{2N}{\pi}$')
```

```
plt.xlabel('Number of Steps (N)')
plt.ylabel('Distance Squared (D^2)')
plt.title('Random Walk Simulation')
plt.legend()
plt.show()
```



[ ]: