

# Probability & Statistics Project

Hana Ali Rashid, hr05940

Tasmiya Malik, tm06183

Ifrah Ilyas, ii06178

April 27, 2021

## Q1: Random Walk

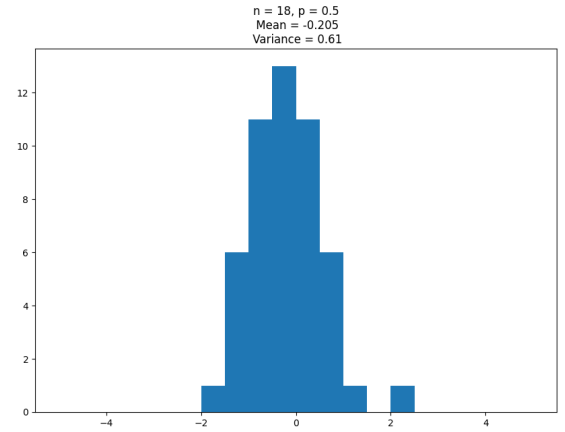
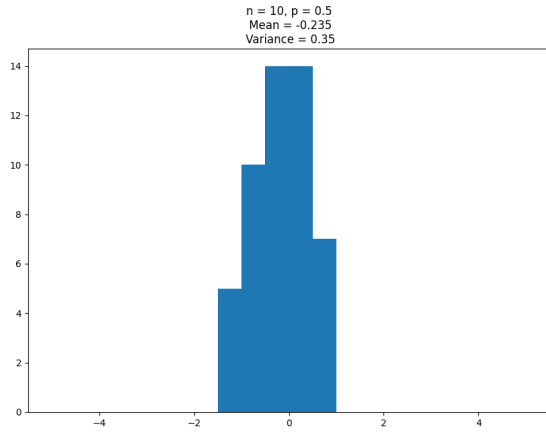
### 1.1

In order to simulate a random walk, the following function takes in values of  $n$  (number of steps to take), and  $p$  (the probability of the object moving one step to the right). It then generates a random number and uses the given probability to determine the direction that the object moves in, and returns the final position.

```
1 def get_updated_position(n,p):
2     pos = 0 #position
3     for _ in range(n):
4         rand = random.randint(1,100) #generating a random number in the range 1 to 100
5         if rand < p*100:
6             pos += 1 #move one step right
7         else:
8             pos -= 1 #move one step left
9     return pos #return final position
```

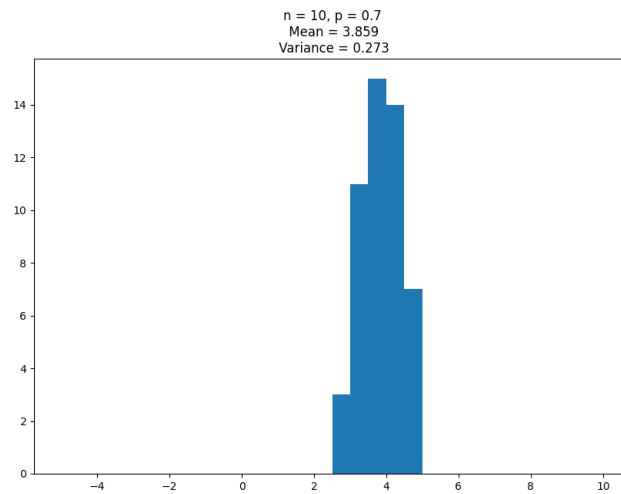
The following function uses the `get_updated_position()` function to get 25 outcomes which are then used to calculate a single expected value for the final position. This is repeated 50 times and the histogram of all these expected values is then plotted.

```
1 def main_11():
2     '''Calculating and plotting expected outcomes for various combinations of n and p'''
3     p = 0.7
4     n = 10
5     expected = []
6     for j in range(50): #expected values for each (n,p) for 50 iterations
7         outcomes = []
8         for i in range(25):
9             outcomes.append(get_updated_position(n,p))
10            expected.append(sum(outcomes)/25) #appending the expected (average) value for each(n,p)
11        #plotting and showing a histogram of calculated expected values
12        fig, ax = plt.subplots(figsize=(10, 7))
13        ax.hist(expected, bins = range(-5,10))
14        plt.title('n = '+str(n)+', p = '+str(p)+'\n Mean = '+str(round(statistics.mean(expected),3))+'\n Variance = '+str(round(statistics.variance(expected),3)))
15        plt.savefig("Q1_histograms/q1"+'_n = '+str(n)+'_ p = '+str(p)+'.png')
16        plt.show()
```



(a) No. of steps taken = 10 with probability of moving a step right = 0.5.

(b) No. of steps taken = 18 with probability of moving a step right = 0.5.



(c) No. of steps taken = 10 with probability of moving a step right = 0.7.

Histograms produced by the above code for expected final positions of objects against frequency, where  $n$  is the number of steps taken and  $p$  is the probability of the object moving one step to the right.

Histograms (a) and (b) have the same value of  $p$  and varying number of steps  $n$ .

Both histograms appear to follow a normal distribution, with (a) having a mean of  $-0.235$  and a variance of  $0.35$  and (b) having a mean of  $-0.205$  and a variance of  $0.61$ . Increasing the number of steps has not had a significant effect on the mean but did increase the variance.

Histograms (a) and (c) have the same number of steps  $n$  and varying values of  $p$ .

Histogram (c) also appears to follow a normal distribution, having a mean of  $3.859$  and a variance of  $0.273$ . Increasing  $p$  appears to have affected the mean of the distribution but not the variance as such.

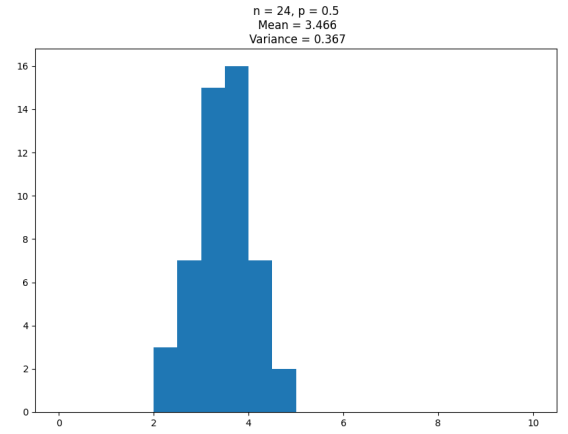
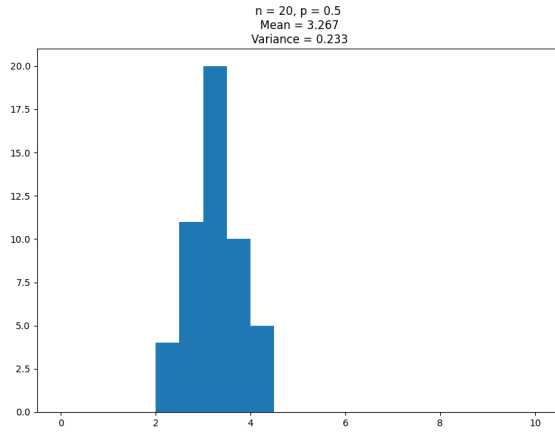
## 1.2

The following function is the modified version of the `get_updated_position()` function from the previous part such that the object cannot move into the negative part of the number line.

```
1 def get_updated_position_restricted(n,p):
2     pos = 0 #position
3     for _ in range(n):
4         rand = random.randint(1,100) #generating a random number in the range 1 to 100
5         if rand < p*100 or pos <= 0: #move one step right if pos == 0
6             pos += 1
7         else:
8             pos -= 1 #move one step left
9     return pos #return final position
```

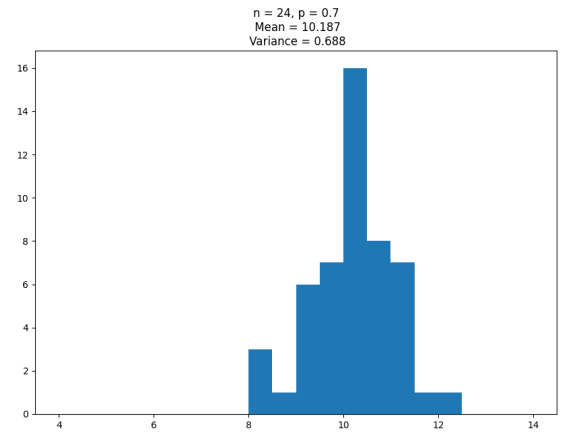
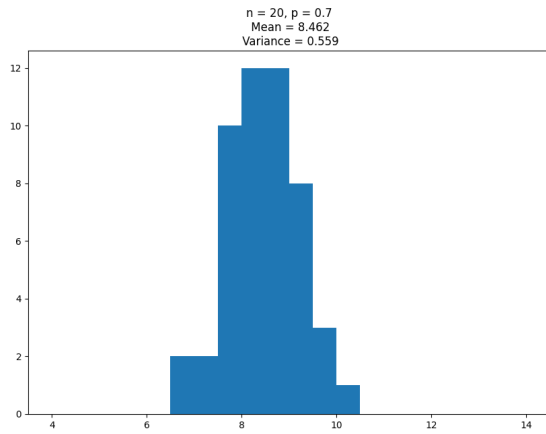
The following function uses the `get_updated_position_restricted()` function to get 25 outcomes which are then used to calculate a single expected value for the final position. This is repeated 50 times and the histogram of all these expected values is then plotted.

```
1 def main_12():
2     '''Calculating and plotting expected outcomes for various combinations of n and p'''
3     p = 0.7
4     n = 24
5     expected = []
6     for j in range(50): #expected value for each (n,p) for 25 iterations
7         outcomes = []
8         for i in range(25):
9             outcomes.append(get_updated_position_restricted(n,p))
10        expected.append(sum(outcomes)/25) #appending the expected (average) value for each(n,p)
11    #plotting and showing a histogram of calculated expected values
12    fig, ax = plt.subplots(figsize =(10, 7))
13    ax.hist(expected, bins =
14    [4,4.5,5,5.5,6,6.5,7,7.5,8,8.5,9,9.5,10,10.5,11,11.5,12,12.5,13,13.5,14])
15    plt.title('n = '+str(n)+', p = '+str(p)+'\n Mean = '+str(round(statistics.mean(expected)
16    ,3))+'\n Variance = '+str(round(statistics.variance(expected),3)))
17    plt.savefig("Q1_histograms/Q1.22 "+"_n = '+str(n)+'_p = '+str(p)+'.png')
18    plt.show()
```



(d) No. of steps taken = 20 with probability of moving a step right = 0.5.

(e) No. of steps taken = 24 with probability of moving a step right = 0.5.



(f) No. of steps taken = 20 with probability of moving a step right = 0.7.

(g) No. of steps taken = 24 with probability of moving a step right = 0.7.

Histograms produced by the above code for expected final positions of objects against frequency, where  $n$  is the number of steps taken and  $p$  is the probability of the object moving one step to the right.

Histograms (d) and (e) have the same value of  $p$  and varying number of steps  $n$ . Both histograms appear to follow a normal distribution, with (d) having a mean of 3.267 and a variance of 0.233 and (e) having a mean of 3.466 and a variance of 0.367. Again, increasing the number of steps has not had a significant effect on the mean but did increase the variance.

Histograms (f) and (g) have the same number of steps  $n$  respectively as (d) and (e) and varying values of  $p$ . They also appear to follow a normal distribution, with (f) having a mean of 8.462 and a variance of 0.559 and (g) having a mean of 10.187 and a variance of 0.668. Increasing  $p$  appears to have increased the mean and variance, but the mean has increased more significantly.

Additionally, due to the added constraint in this part, we see that no expected value for the final position is negative.

### 1.3

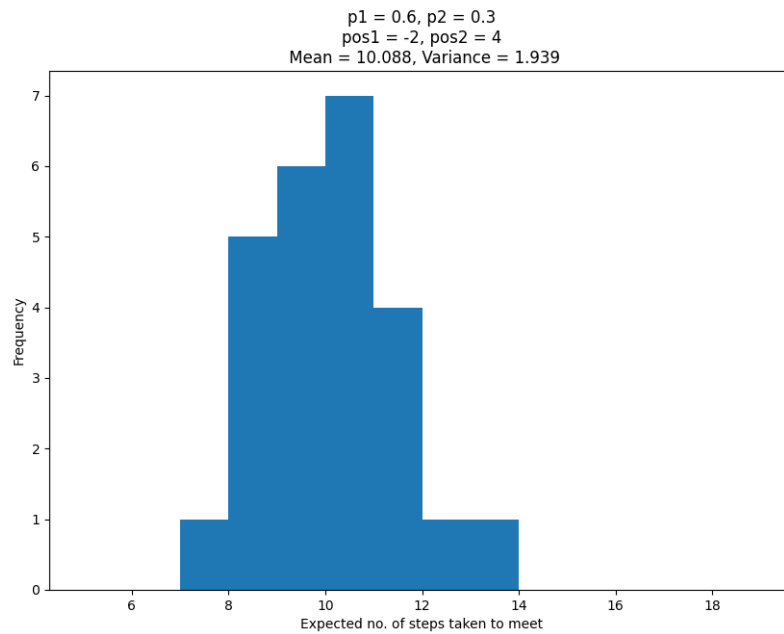
In this part, we are required to find the number of steps taken by two objects initialized at different points, and moving randomly to meet. The following function takes in the initial position of both objects *pos1* and *pos2* as well as their probabilities *p1* and *p2*, generates random numbers and determines the number of steps taken to meet.

```
1 def stepsToMeet(pos1,pos2,p1,p2):
2     count = 0 #keeps count of number of steps taken for objects to meet
3     while pos1 != pos2:
4         rand = random.randint(1,100) #generating a random number in the range 1 to 100 to
         determine outcome
5         if rand < p1*100:
6             pos1 += 1 #move one step right
7         else:
8             pos1 -= 1 #move one step left
9         rand = random.randint(1,100) #generating a random number in the range 1 to 100
10        if rand < p2*100:
11            pos2 += 1 #move one step right
12        else:
13            pos2 -= 1 #move one step left
14        count += 1
15    return count
```

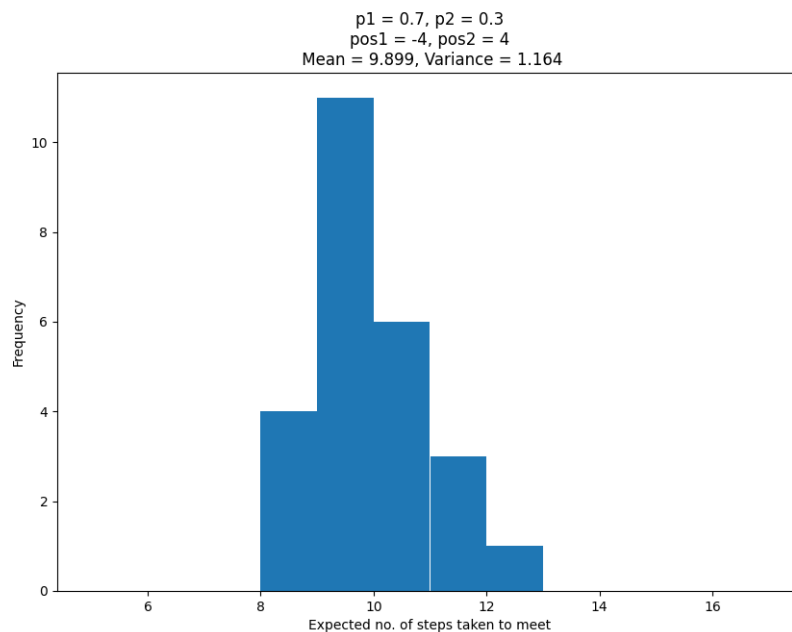
The following function uses the `stepsToMeet` function and calculates 25 expected values using 25 outcomes based on the values of *pos1*, *pos2*, *p1* and *p2*. It then plots a histogram of these expected values as well.

```
1 def main_13():
2     '''Calculating and plotting expected outcomes for various combinations of n & p'''
3     expected = []
4     p1 = 0.7
5     p2 = 0.3
6     pos1 = -4
7     pos2 = 4
8     for i in range(25): #calculating the expected value for each (n,p) for 25 iterations
9         outcomes = []
10        for j in range(25):
11            outcomes.append(stepsToMeet(pos1,pos2,p1,p2))
12        #calculating the average expected value for each(n,p)
13        expected.append(sum(outcomes)/25) #appending the expected (average) value for each(n,p)
14    #plotting a histogram of calculated expected values
15    fig, ax = plt.subplots(figsize =(10, 7))
16    ax.set_xlabel('Expected no. of steps taken to meet')
17    ax.set_ylabel('Frequency')
18    ax.hist(expected, bins = range(5,18))
19    plt.title('p1 = '+str(p1)+' , p2 = '+str(p2)+'\npos1 = '+str(pos1)+' , pos2 = '+str(pos2)+'\n
        Mean = '+str(round(statistics.mean(expected),3))+' , Variance = '+str(round(statistics.
        variance(expected),3))')
20    plt.savefig("Q1_histograms/Q1.3 "+'_p1 = '+str(p1)+'_p2 = '+str(p2)+'_pos1 = '+str(pos1)+'
        _pos2 = '+str(pos2)+'(6).png')
21    plt.show()
```

Histograms produced by the above code for expected number of steps taken for two objects walking randomly to meet, against frequency:



The above histogram shows that the expected number of steps to meet for two objects follows a normal distribution of mean 10.09 and variance 1.94 when they begin 6 steps apart with the given probabilities.

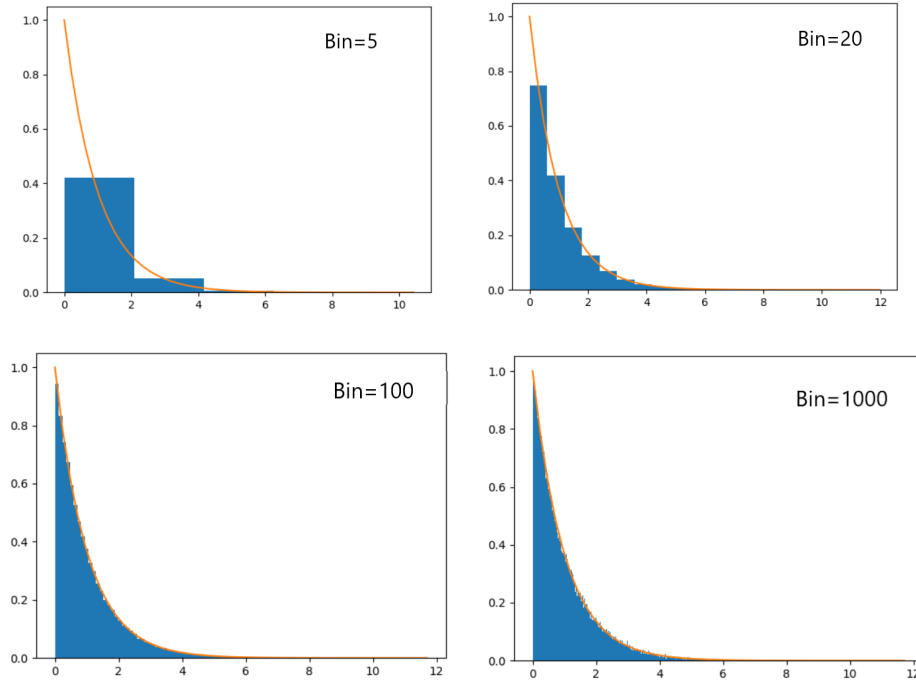


The above histogram shows that the expected number of steps to meet for two objects follows a normal distribution of mean 9.90 and variance 1.16 when they begin 8 steps apart with the given probabilities.

## Q2: Simulating Distributions

### 2.1

It does accomplish as it plots for a value of  $x$  ( $y$  is the variable expressed in terms of  $x$ ) against the pdf  $e^{-y}$  such that each bar in the histogram shows the probability for an average value of  $x$ . The distribution plotted is a continuous exponential random variable. The number of bins used are significant for our simulation which will become evident with the following images:



**Figure 1:** Simulations done with different number of bins

It is noticeable that, with more number of bins, the accuracy of the results increases. However, it is also interesting to note that using 1000 bins almost gives the same result as 100 bins. Therefore, while simulating we need to take an appropriate values for bins for better approximation and less running time.

### 2.2

In the code we see that :

$$y = \frac{1}{1-x}$$

With the trick shown, we can back track to find the original distribution.

$$1 - \frac{1}{y} = x$$

We know that the upper bound will be  $y$ . For the lower bound we see that:

$$\begin{aligned} \frac{1}{1} - \frac{1}{y} &= x \\ -\left(\frac{1}{y} - \frac{1}{1}\right) &= x \\ \therefore -\frac{1}{y} \Big|_1^y &= x \end{aligned}$$

Performing derivative on the left hand side:

$$\frac{d}{dy} \frac{-1}{y} = (-1) \frac{-1}{y^2} = \frac{1}{y^2}$$

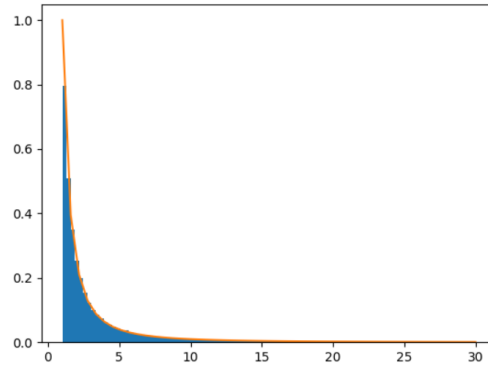
$$\therefore \int_1^y \frac{1}{y^2} dy = x$$

We can affirm that we have reached the right deduction by looking at line 12 of the code. So, the distribution of y is :

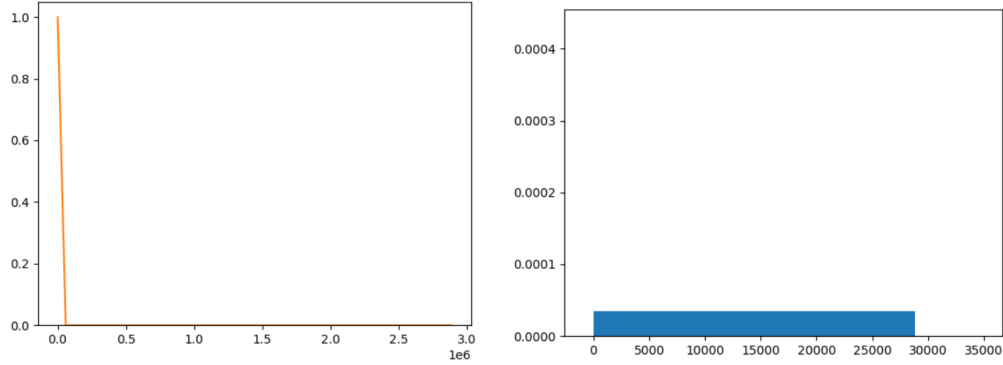
$$f_Y(y) = \frac{1}{y^2} \text{ for } y \geq 1$$

```
1 # y . sort ()
2 # ind = ( np . array ( y ) > 30 ) . tolist () . index (1)
3 # y = y [: ind ]
```

Lines 5-7 is simply trying to implement a cutoff for values of y greater 30 as values starting to get more smaller. Given that we taken small amount of bins, the average value of each bar becomes negligible and is not visible in the graph. Therefore we cut it down to 30 so that the approximation is more accurate.



**Figure 2:** Before Removing



**Figure 3:** After Removing

## 2.3

For,

$$f_Y(y) = \frac{1}{y^3} \text{ for } y \geq \sqrt{\frac{1}{2}}$$

$$P(Y < y) = P(X < x)$$

$$\int_{\sqrt{\frac{1}{2}}}^y \frac{1}{y^3} dy = x$$

$$\left. \frac{-1}{2y^2} \right|_{\sqrt{\frac{1}{2}}}^y = x$$

$$\frac{-1}{2y^2} + \frac{1}{2(\sqrt{\frac{1}{2}})^2} = x$$

$$\frac{-1}{2y^2} + 1 = x$$



$$\frac{1}{2y^2} = 1 - x$$

$$y^2 = \frac{1}{2(1-x)}$$

$$y = \sqrt{\frac{1}{2(1-x)}}$$

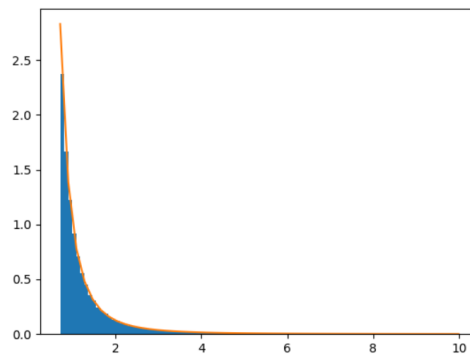
Code:

```

1 y = []
2 for i in range (100000) :
3     x = np . random . random ()
4     y . append ((1 / (2*(1 - x )))**(1/2) )
5
6 y . sort ()
7 ind = ( np . array ( y ) > 10 ) . tolist () . index (1)
8 y = y [: ind ]
9 bins = 100
10 binWidth = (max( y ) - min( y ) ) / bins
11 plt . hist ( y , bins = bins , weights = np . ones ( len ( y ) ) / ( len ( y ) * binWidth ) )
12 values = np . linspace ( min( y ) , max( y ) , 50)
13 plt . plot ( values , (1 / values ** 3) )
14 plt . show ()

```

The output of simulation for this y:



**Figure 4:** Simulation of distribution  $f_Y(y) = \frac{1}{y^3}$  for  $y \geq \sqrt{\frac{1}{2}}$

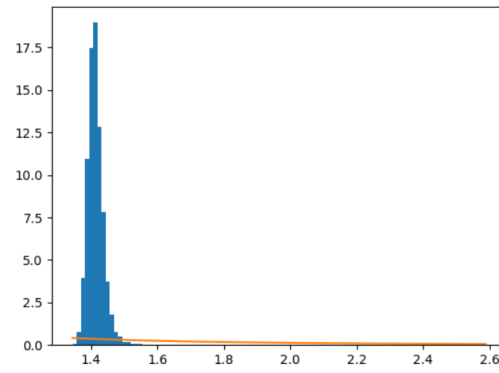
For the simulation of expected values, the same code used for the above function ran for 10000 times and each run would take the mean of y :

```

1 y = []
2 for j in range (10000) :
3     x = np . random . random (10000)
4     x=(1 / (2*(1 - x )))**(1/2)
5     y.append(np.mean(x))
6 print(y)
7 bins = 100
8 binWidth = (max( y ) - min( y ) ) / bins
9 plt . hist ( y , bins = bins , weights = np . ones ( len ( y ) ) / ( len ( y ) * binWidth ) )
10 values = np . linspace ( min( y ) , max( y ) , 50)
11 plt . plot ( values , (1 / values ** 3) )
12 plt . show ()

```

The output of this simulation:



**Figure 5:** Simulation of expected values of  $fY(y) = \frac{1}{y^3} f \text{or } y \geq \sqrt{\frac{1}{2}}$

## Q3. Picking a Random Point Correctly

### 3.1

For part 3.1, we followed the method specified in the question. The program takes radius as an input from the user, and passes it to *gen\_points(R)* to display the graph. It first finds random radius  $R$  and angle  $\theta$ , then converts it into cartesian coordinates. It stores the values of  $x$  and  $y$  till the iterations are running, then plot the points. To make sure the points lie in the given radius, we also draw the circle.

- *polar\_to\_cart(r,t)*: This function takes the polar coordinates ( $r$  and  $\theta$ ) and converts them to their cartesian equivalent. It is called from *random\_point(R)*, to return the cartesian coordinates. It uses the following formulae,

$$x = r \cos(\theta)$$

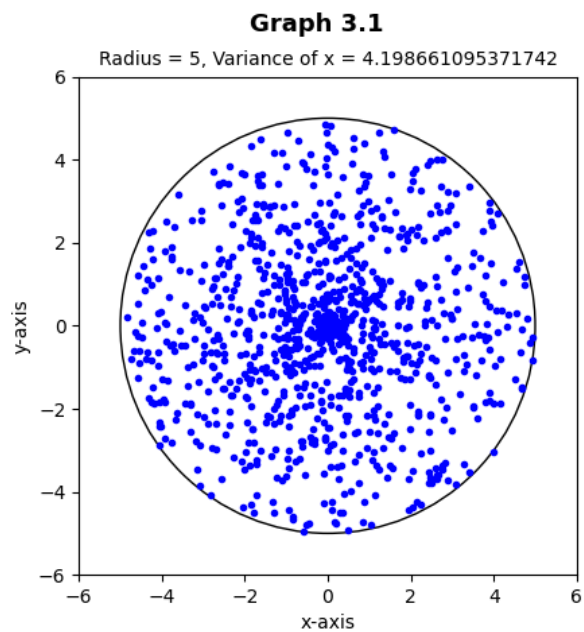
$$y = r \sin(\theta)$$

- *random\_point(R)*: The function is responsible for generating a random point in polar coordinates, and returns the converted cartesian coordinates.
- *gen\_points(R)*: This is the main function which runs 1000 iterations, and calls *random\_point(R)*, and stores the result in the separate lists for  $X$  and  $Y$  coordinates. Then the lists are passed to the plot functions, to graph the points.

The resulting graph has very dense distribution towards the center.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import math as m
4
5 def random_point(R):
6
7     #R=R, theta = angle
8     r = np.random.uniform(0, R) #random R in range 0-R
9     theta = np.random.uniform(0, 360)*(m.pi/180) #random theta between 0-360 degrees
10    return(polar_to_cart(r,theta))
11
12 def polar_to_cart(r,t):
13
14     x = r * np.cos(t) #x = rcos(theta)
15     y = r * np.sin(t) #y = rsin(theta)
16     return((x,y))
17
18 def gen_points(R):
19
20     I = 1000 #iterations
21     X = [] #list for x axis
22     Y = [] #list for y axis
23
24     #generating the lists for X and Y
25     for _ in range(I):
26         point = random_point(R)
27         X.append(point[0])
28         Y.append(point[1])
29
30     """Drawing the circle and plotting the points"""
31
32     #Drawing the circle
33     fig, ax = plt.subplots()
34     ax.axis([-R-1,R+1,-R-1,R+1])
35     circle = plt.Circle((0,0),R,Fill=False)
36     ax.set_aspect(1)
37     ax.add_artist(circle)
38     plt.title("circle")
39
40     #Plotting X and Y
41     plt.plot(X,Y,'.',color="blue")
42     fig.suptitle("Graph 3.1",fontsize=13,fontweight='bold')
43     plt.xlabel("x-axis")
```

```
44 plt.ylabel("y-axis")
45 plt.title("Radius = {}, Variance of x = {}".format(R,np.var(X)),fontsize=10)
46 plt.savefig("Q3/Q3(1).png")
47 plt.show()
```



## 3.2

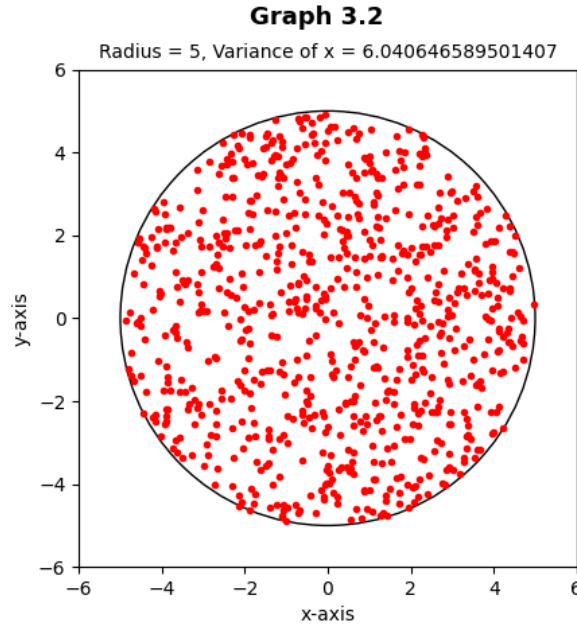
In part 3.2, we calculated the cartesian coordinates directly from the specified radius  $R$ . If the random point lied farther from the distance from the origin, it will be discarded and the program will find a new point.

- *random\_point(R)*: This function finds random points for  $x$  and  $y$  coordinates between the range  $-R$  and  $R$ . It returns the tuple with cartesian coordinates.
- *dist\_from\_origin(x,y)*: The function find and return the distance of the point from the origin.

$$dist = \sqrt{x^2 + y^2}$$

- *gen\_points(R)*: The main function iterates 1000 times and calls *random\_point(R)*. Then it checks the condition that if the distance of the point  $(x,y)$  is within the radius  $R$ , it would append the points to their respective lists, else the counter will be decremented by 1 and the point will not be added to the list. Then the lists  $X$  and  $Y$  are passed to the plot functions for display.

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 import math as m
4
5 def random_point(R):
6
7     x = np.random.uniform(-R,R) #random point x
8     y = np.random.uniform(-R,R) #random point y
9     return((x,y))
10
11 def dist_from_origin(x,y):
12     return(m.sqrt(x**2+y**2))
13
14 def gen_points(R):
15
16     X = []          #list for x-axis
17     Y = []          #list for y-axis
18     I = 1000        #iterations
19
20     for a in range(I):
21         point = random_point(R)
22         if dist_from_origin(point[0],point[1]) <= R: #if distance is smaller than or equal to
23             R
24             X.append(point[0])
25             Y.append(point[1])
26         else: #if distance is greater than R
27             a = a-1
28
29     """Drawing the circle and plotting the points"""
30
31     #Drawing the circle
32     fig, ax = plt.subplots()
33     ax.axis([-R-1,R+1,-R-1,R+1])
34     circle = plt.Circle((0,0),R,Fill=False)
35     ax.set_aspect(1)
36     ax.add_artist(circle)
37     plt.title("circle")
38
39     #Plotting the points
40     plt.plot(X,Y,'.',color="red")
41     fig.suptitle("Graph 3.2",fontsize=13,fontweight='bold')
42     plt.xlabel("x-axis")
43     plt.ylabel("y-axis")
44     plt.title("Radius = {}, Variance of x = {}".format(R,np.var(X)),fontsize=10)
45     plt.savefig("Q3/Q3(2).png")
46     plt.show()
```



### 3.3

In part 3.3, we had to modify the function for polar coordinates, such that the graph it made was uniformly distributed, similar to the histogram of part 3.2.

The reason behind the uneven distribution when we used polar coordinates was because of the way polar and cartesian planes interact. If we plot points in a smaller radius with varying angles, the points will be closer, as compare to a larger radius. For larger radius the points will be far apart. So when we were converting them into cartesian coordinates, many of them lied near the origin. To solve this, we simply take the CDF, i.e. the area of the circle, and make radius  $r$  its subject. Mathematically,

$$P(r \leq R) = k.\pi.r^2$$

To find  $k$ , we use  $P(r=R)$ , which is equals to 1.

$$k.\pi.R^2 = 1$$

$$k = \frac{1}{R^2\pi}$$

Now that we know the value of  $k$ , we can find the whole CDF by plugging its value.

$$CDF = \left(\frac{1}{R^2\pi}\right)(r^2\pi) = \frac{r^2}{R^2}$$

Now we make radius  $r$ , the subject of the equation.

$$r = \sqrt{R^2.CDF} = R\sqrt{CDF}$$

We can find our random radius  $r$  using this formula. This will produce the points in proportionality with the area of the circle, thus keeping the point uniform.

From here, we will describe the functionality of the coded program.

- `random_point(R)`: The function finds random radius  $R$  and theta  $\theta$ . It uses the function for  $R$  derived above.

- *polar\_to\_cart(cord)*: This is the same function that we used in part 3.1. The function takes radius  $r$  and angle  $t$  as the input, and returns the cartesian coordinates. We use the following formulae to derive the coordinates,

$$x = r \cos(t)$$

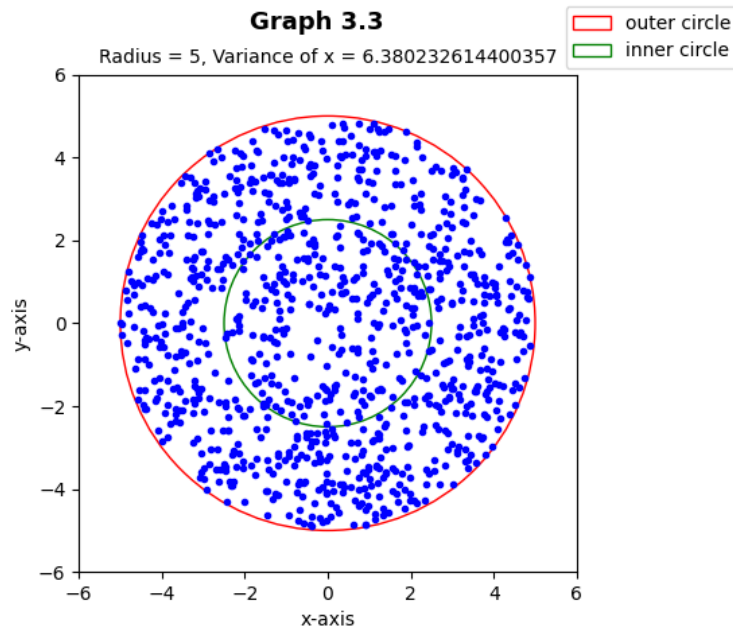
$$y = r \sin(t)$$

- *gen\_points(R)*: Like previous parts, this function makes 1000 iterations, and call *random\_point(R)* for each of them, then stores the returned values in their respective lists X and Y. Then it draws outer and inner circles of the given radius  $R$  and  $\frac{R}{2}$  respectively, and then plot the points stored in the lists X and Y.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import math as m
4
5 def random_point(R):
6
7     #r=radius, theta=angle
8     r = R*m.sqrt(np.random.uniform(0,1))    #random r using R(CDF)**1/2
9     theta = np.random.uniform(0, 360)*(m.pi/180)    #random theta between 0-360 degrees
10    return(polar_to_cart(r,theta))
11
12 def polar_to_cart(r,t):
13
14    x = r * np.cos(t)    #x = rcos(theta)
15    y = r * np.sin(t)    #y = rsin(theta)
16    return((x,y))
17
18 def gen_points(R):
19
20    I = 1000    #iterations
21    X = []    #list for x axis
22    Y = []    #list for y axis
23
24    #generating the lists for X and Y
25    for _ in range(I):
26        point = random_point(R)
27        X.append(point[0])
28        Y.append(point[1])
29
30    """Drawing the circles and plotting the points"""
31
32    #Outer Circle of radius R
33    fig, axes = plt.subplots()
34    axes.axis([-R-1,R+1,-R-1,R+1])
35    circle1 = plt.Circle((0,0),R,Fill=False,color="red")
36    axes.set_aspect(1)
37    axes.add_artist(circle1)
38
39    #inner circle of radius R/2
40    circle2 = plt.Circle((0,0),R/2,Fill=False,color="green")
41    axes.set_aspect(1)
42    axes.add_artist(circle2)
43    fig.legend([circle1,circle2],["outer circle", "inner circle"])
44
45    #Plotting X and Y
46    plt.plot(X,Y,'.',color="blue")
47    fig.suptitle("Graph 3.3", fontsize=13,fontweight='bold')
48    plt.xlabel("x-axis")
49    plt.ylabel("y-axis")
50    plt.title("Radius = {}, Variance of x = {}".format(R, np.var(X)),fontsize=10)
51    plt.savefig("Q3/Q3(3).png")
52    plt.show()

```



## Q4. Saying Random is not enough

### 4.1

For 4.1, we took the approach specified in the question. We first pick two random angles  $\theta_1$  and  $\theta_2$ . Then find their difference to find the angle lying in between  $\theta_1$  and  $\theta_2$ . Then we use this new angle in *cord length* formula. The program does this process 1000 times and stores the values in the specified list, then graphs a histogram for the length of the cords and the probability of them occurring. Following describes each function and their purpose/method for the program.

1. *random\_theta()*: The function finds random  $\theta_1$  and  $\theta_2$ .  $\theta$  is chosen randomly from  $\text{degrees}(\phi)$  and then converted to radians.

$$\theta_1 = \phi_1 \cdot \frac{\pi}{180}$$

$$\theta_2 = \phi_2 \cdot \frac{\pi}{180}$$

2. *cord(R)*: The function takes the radius as the input and returns length of the cord between  $\theta_1$  and  $\theta_2$ . It first calls *random\_theta()*, and find the absolute value of their difference ( $\theta = |\theta_1 - \theta_2|$ ). Then the length of the cord is calculated using the following formula for cord length,

$$\text{cord length} = 2R \cdot \sin(\theta/2)$$

3. *find\_cords1(R)*: The function takes radius R as input, then draw the histogram as the output. It iterates 1000 times, calls *cord(R)* for each iteration and append the result in *cord\_len*. The bin value is set to be 50 as it provides a good threshold for the 1000 values. Then it plots the graph using matplotlib.

```

1
2 # 4.1
3
4 def random_theta():
5     theta1 = np.random.uniform(0,360)*(m.pi/180)    #random theta 1
6     theta2 = np.random.uniform(0,360)*(m.pi/180)    #random theta 2
7     return (theta1,theta2)
8
9 def cord(R):
10     angle = random_theta()
11     theta = abs(angle[0] - angle[1])    #theta between the two radii(theta1 and theta2)
12     l = 2*R*m.sin(theta/2)    #length of cord = 2rsin(theta/2)
13     return(l)

```

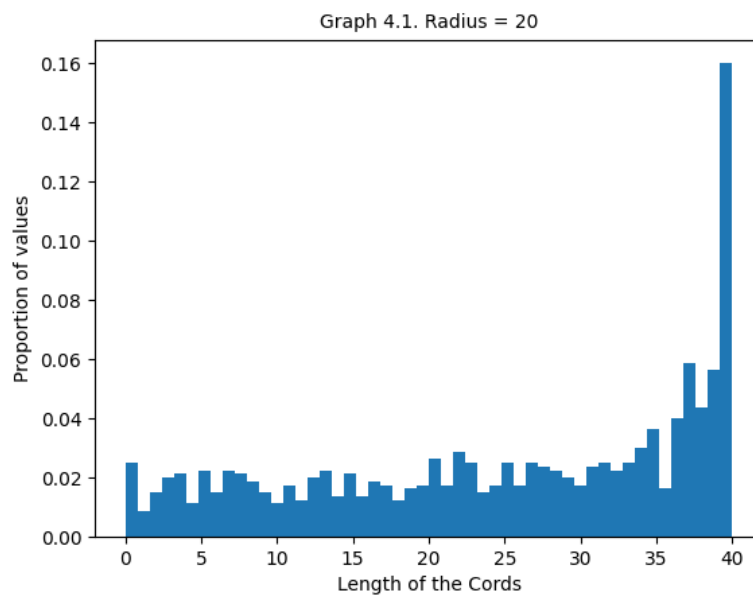


```

13
14 def find_cords1(R):
15
16     cord_len = []
17     I = 1000    #iterations
18
19     for _ in range(I):
20         cord_len.append(cord(R))
21
22     #Finding bins
23     bin_val = 50
24
25     plt.hist(cord_len, bins = bin_val,density=True)
26     plt.title("Graph 4.1. Radius = {}".format(R),fontsize=10)
27     plt.ylabel("Proportion of values")
28     plt.xlabel("Length of the Cords")
29     plt.savefig("Q4/Q4(1).png")
30     plt.show()

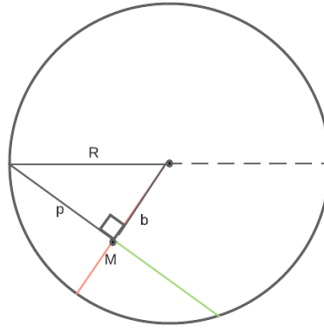
```

**Figure 6:** Histogram of 4.1



## 4.2

For 4.2, we follow the instructions given in the question. We first pick a direction/angle and assume an imaginary radius there. Then we pick a random point on that radius and find its distance from the origin (0,0). This serves as our base, while to radius R serves as our hypotaneous. Then we find the perpendicular using the pythagorean theorem, and returns the twice of it. As the perpendicular shows the half part of that cord, with randomly chosen point as its center. This point is also the midpoint for our circle, so we double the perpendicular to find the whole cord length.



1. *random\_cord(R)*: This function does multiple tasks. It first find a random direction/angle in the circle, assumes a radius on that angle, and pick a random point on the assumed radius. This point serves as the Midpoint M as well. Then it converts the polar coordinates into cartesian using the following formulae,

$$x = r\cos(\theta)$$

$$y = r\sin(\theta)$$

Then we find our base b and perpendicular p, while assuming it is a right-angle triangle because *p* and *b* makes a right-angle triangle. In our case we assume the distance between the origin and the random point (*x*, *y*) to be our base, and radius R as our hypotaneous. We use the *distance formula* to find the value of base,

$$base = \sqrt{x^2 + y^2}$$

Then we find the value of perpendicular line by using the pythagorean theroem,

$$perpendicular = \sqrt{R^2 - base^2}$$

Then we simply multiply the value of perpendicular to find the length of the cord and return it.

2. *find\_cords2(R)*: This function runs 1000 iterations of the *random\_cord(R)* function, and stores their result in the list *cord\_len*, sets the bin width as 50, and then graphs the histogram as the result.

We run this program with radius 30. The resulting histogram seems to be an exponentially rising graph, with a flat start.

```

1
2 # 4.2
3
4
5 def random_cord(R):
6     theta = np.random.uniform(0,360)*(m.pi/180)
7     point_at_radius = np.random.uniform(0,R) #point at R
8
9     #cartesian cordinates at the picked point
10    x = point_at_radius*m.cos(theta)
11    y = point_at_radius*m.sin(theta)
12
13    #finding base line from center to point_at_radius using distance formula
14    base = m.sqrt(x**2+y**2)

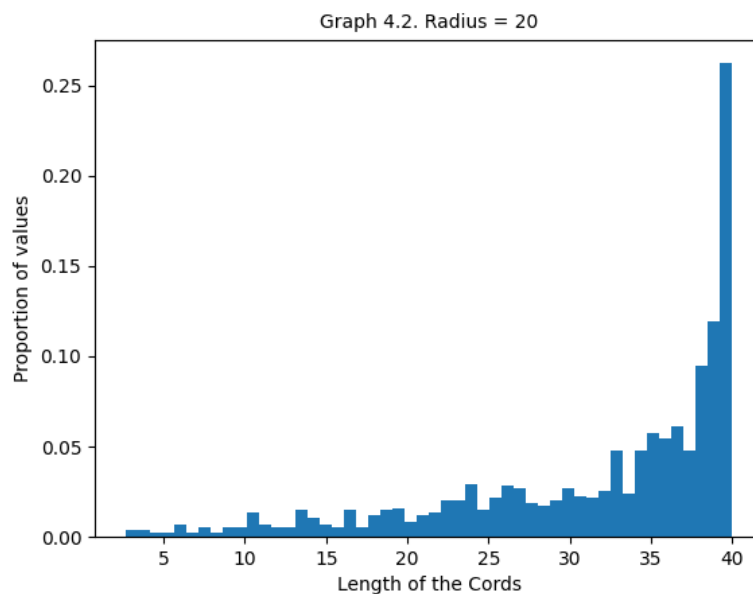
```

```

14
15 #perpendicular using pythagorean theorem, p = sqrt(h^2-b^2)
16 perp = m.sqrt(R**2-base**2)
17
18 #length of the cord
19 l = 2*perp
20
21 return l
22
23 def find_cords2(R):
24
25     cord_len = []
26     I = 1000    #iterations
27
28     for _ in range(I):
29         cord_len.append(random_cord(R))
30
31     #finding the bin size
32     bin_val = 50
33
34     plt.hist(cord_len, bins = bin_val,density=True)
35     plt.title("Graph 4.2. Radius = {}".format(R),fontsize=10)
36     plt.ylabel("Proportion of values")
37     plt.xlabel("Length of the Cords")
38     plt.savefig("Q4/Q4(2).png")
39     plt.show()

```

**Figure 7:** Histogram of 4.2



## 4.3

For 4.3, it is similar to 4.2 with a slight difference. Instead of determining a direction, we pick a random point in the circle, and then calculates its distance from the origin (0,0). If the distance is within the circle as well, we regard it as our base/adjacent of the right triangle. Then we move towards calculating the opposite/perpendicular by using pythagorean theorem, multiply it by 2 to find the total length, and then store it in the list. Then we pass that list plotting functions to plot a histogram.

1. *p\_to\_o(cord)*: This function simply returns the distance between the randomly chosen point and the origin (0,0). It takes a tuple of cartesian coordinates as the input and returns the distance. It uses the distance formula,

$$dist = \sqrt{x^2 + y^2}$$

2. *random\_point(R)*: It chooses random points for cartesian coordinates (x,y) between the R and -R.
3. *cal\_cord(R,pnt)*: It takes the randomly chosen point, and finds its distance from the origin, then uses it as the adjacent/base of our right-angle triangle within the circle. Then it finds the value of the opposite/perpendicular via pythagorean theorem, and returns the twice of the opposite as the cord length.
4. *find\_cords3(R)*: This function runs 1000 times while calling the *random\_point(R)* function, then it checks if the distance/base lies within the circle or not. If it does, then it proceeds to find *cal\_cord(R,pnt)* and store it, else it decrements the counter by 1, to discard that value and find another.

We run this program for radius 20. The histogram looks quite linear.

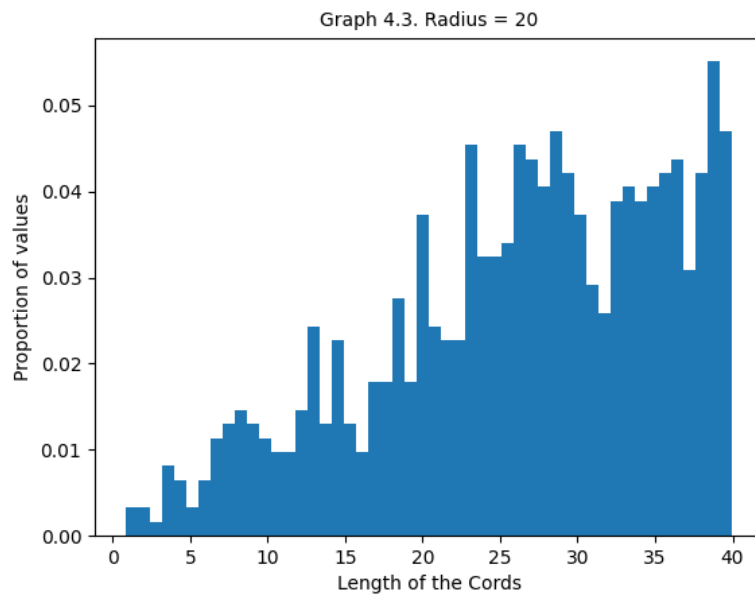
```
1
2 # 4.3
3 -----
4 def p_to_o(cord):
5     return m.sqrt(cord[0]**2+cord[1]**2)
6
7 def random_point(R):
8
9     x = np.random.uniform(-R,R) #random point x
10    y = np.random.uniform(-R,R) #random point y
11
12    return (x,y)
13
14 def cal_cord(R,pnt):
15     #finding adjacent from the random point.
16     adj = p_to_o(pnt)
17
18     #length of opposite
19     opp = m.sqrt(R**2-adj**2)
20
21     #length of the cord
22     l = 2*opp
23
24     return l
25
26 def find_cords3(R):
27
28     I = 1000
29     cord_len = []
30
31     for a in range(I):
32         point = random_point(R)
33         if p_to_o(point) <= R:
34             cord_len.append(cal_cord(R,point))
35         else:
36             a = a - 1
37
38     #finding the bin sizes
39     bin_val = 50
40
41     plt.hist(cord_len, bins = bin_val,density=True)
42     plt.title("Graph 4.3. Radius = {}".format(R),fontsize=10)
```

```

43 plt.ylabel("Proportion of values")
44 plt.xlabel("Length of the Cords")
45 plt.savefig("Q4/Q4(3).png")
46 plt.show()

```

**Figure 8:** Histogram of 4.3



## 4.4

Out of all the distributions, I think the third approach (part 4.3) is the best to take. Its graph seems somewhat linear, which will ensure a little uniformity. As the previous graphs showed dense distribution towards the center of the circle, i.e. near the diameter, however in 4.3, the graph seems more distributed.

## Q5: Hypothesis Testing

### 5.1

This question requires us to use hypothesis testing to determine whether the null hypothesis that a coin is fair is true or not.

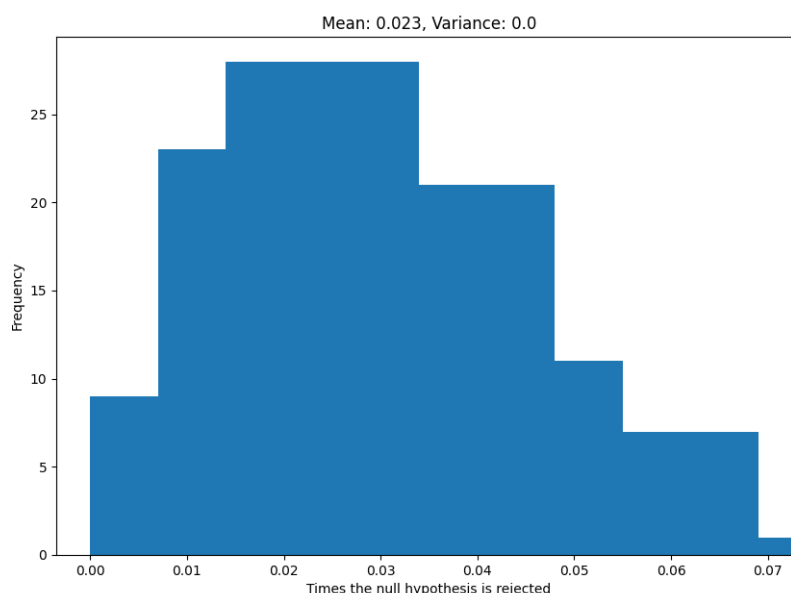
The following function implements the simulation of a fair coin:

```
1 def cointoss():
2     #assuming 0 = T & 1 = H
3     return(random.randint(0,1))
```

This function then uses the above function to simulate 10 coin tosses 100 times and finds the expected number of times the null hypothesis is rejected for each iteration:

```
1 def simulate_tosses():
2     expected = []
3     for n in range(100): #to get 100 expected values
4         rejected = []
5         for j in range(100): #to get 1 expected value
6             outcomes = []
7             for i in range(10): #to get 10 outcomes/1 rejected value
8                 outcomes.append(cointoss())
9             #check if hypothesis is accepted or rejected
10            if sum(outcomes) == 0 or sum(outcomes) == 1 or sum(outcomes) == 10 or sum(outcomes) == 9:
11                rejected.append(1) #hypothesis rejected
12            else:
13                rejected.append(0) #hypothesis accepted
14            expected.append(sum(rejected)/100)
15        #plotting histogram
16        fig, ax = plt.subplots(figsize=(10, 7))
17        ax.hist(expected,width = 0.01)
18        ax.set_xlabel('Times the null hypothesis is rejected')
19        ax.set_ylabel('Frequency')
20        plt.title('Mean: '+str(round(statistics.mean(expected),3))+', Variance: '+str(round(statistics.variance(expected),3)))
21        plt.savefig("Q5_histograms/Q5.1.png")
22        plt.show()
```

Histogram of expected values:



Probability we will reject the null hypothesis even though it is true:

*Simulation-wise:* According to the mean of the expected values, the probability is 2.3%.

*Mathematically:* We reject the null hypothesis if the probability of the outcomes is below the threshold. Since the probability of the outcomes being accepted is 0.95 and the probability of them being rejected is 0.05, the probability that we reject the null hypothesis even though it is true is also 0.05 which is equal to the threshold.

### 5.2.1

This question requires us to use hypothesis testing to determine the validity of the null hypothesis that the mean length of fish in a lake is 23.<sup>1</sup>

The following function conducts a single hypothesis test by taking a sample of 30 fish and calculating the mean and standard deviation of their lengths to check the following condition to accept or reject the null hypothesis:

$$P(|S - u_0| \geq a) < 0.05$$

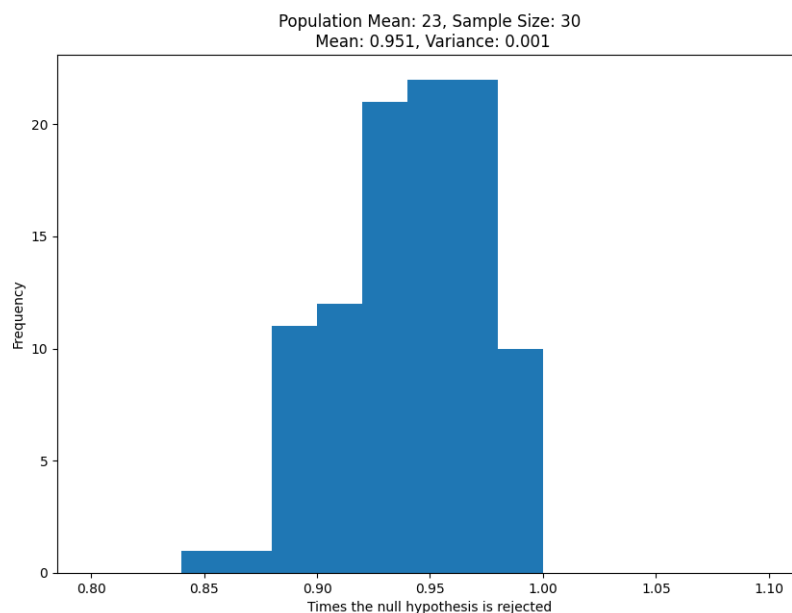
The function returns 1 if the hypothesis is rejected, and 0 if it is accepted.

```
1 def hypothesis_test(u_0, n):
2     s = []
3     for i in range(n): #catching a sample of 30 fish
4         s.append(f.fish())
5     s = np.array(s)
6     std = s.std() #sample standard deviation
7     mean = s.mean() #sample mean
8     a = abs(mean - u_0)
9     prob = 2*(stats.norm(u_0, std/math.sqrt(n)).cdf(u_0 - a))
10    if prob < 0.05:
11        return 1 #rejected
12    else:
13        return 0 #accepted
```

The following function performs 50 experiments and uses the average of their outcome to calculate a single expected value. It repeats this for 100 iterations and plots a histogram of the expected values.

```
1 def experiments():
2     u_0 = 23 #population mean
3     n = 30 #sample size
4     expected = []
5     for j in range(100): #getting 100 expected values
6         outcomes = []
7         for i in range(50): #conducting 50 hypothesis tests
8             outcomes.append(hypothesis_test(u_0,n))
9         expected.append(sum(outcomes)/50) #calculating one expected value from the 50 tests
10    #plotting histogram of expected values
11    fig, ax = plt.subplots(figsize =(10, 7))
12    ax.hist(expected, width = 0.04)
13    ax.set_xlabel('Times the null hypothesis is rejected')
14    ax.set_ylabel('Frequency')
15    plt.title('Population Mean: '+str(u_0)+' , Sample Size: '+str(n)+'\n Mean: '+str(round(
16    statistics.mean(expected),3))+', Variance: '+str(round(statistics.variance(expected),3)))
17    plt.savefig("Q5_histograms/Q5.2.1.png")
18    plt.show()
```

Histogram of expected values of the null hypothesis being rejected:



<sup>1</sup>The fish.pyc file we received was causing errors so we used the new fish.py file provided by a TA



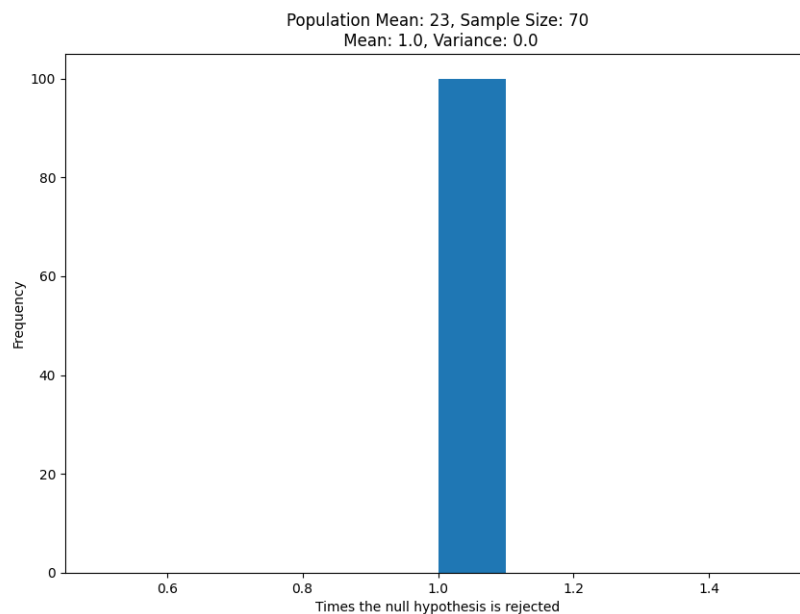
The histogram follows a normal distribution with mean 0.951 and variance 0.001. According to this, the null hypothesis is rejected 95% of the times so it is false. A single hypothesis test may have been sufficient to accept or reject the null hypothesis as the variance is very low.

### 5.2.2

The following function calls the `hypothesis_test()` function defined in 5.2.1 using the same value for the population mean  $u_0$  (23) but passing a value of 70 for the sample size  $n$ . It then performs 50 experiments and uses the average of their outcome to calculate a single expected value. It repeats this for 100 iterations and plots a histogram of the expected values.

```
1 def experiments_updated():
2     u_0 = 23 #population mean
3     n = 70 #sample size
4     expected = []
5     for j in range(100):
6         outcomes = []
7         for i in range(50):
8             outcomes.append(hypothesis_test(u_0,n))
9         expected.append(sum(outcomes)/50)
10    #plotting histogram
11    fig, ax = plt.subplots(figsize =(10, 7))
12    ax.hist(expected, width = 0.1)
13    ax.set_xlabel('Times the null hypothesis is rejected')
14    ax.set_ylabel('Frequency')
15    plt.title('Population Mean: '+str(u_0)+' , Sample Size: '+str(n)+'\n Mean: '+str(round(
16    statistics.mean(expected),3))+', Variance: '+str(round(statistics.variance(expected),3)))
17    plt.savefig("Q5_histograms/Q5.2.2.png")
18    plt.show()
```

Histogram of expected values of the null hypothesis being rejected:



Increasing the value of  $n$  made the variance become 0 where it was previously 0.001. In other words, the null hypothesis is rejected 100% of the times. A single hypothesis test may have been sufficient to reject the null hypothesis in this case as there is no variance in the outcome.

### 5.2.3

**Determining the least value of  $n$  to ensure that the null hypothesis is not wrongly rejected more than 10 percent of the time using simulations.**

The following is a function that returns the length of a fish following a normal distribution with mean 23 and standard deviation 3.

```
1 def my_fish():
2     return np.random.normal(23,3)
```

In order to calculate the expected number of times the null hypothesis is rejected, we use the following function which takes in the value of the population mean  $u_0$  and sample size  $n$  and uses the `my_fish()` function defined above to create the sample for the test. It then conducts the test itself the same way as the previous parts but with a threshold of 0.1.

```

1 def testing(u_0,n):
2     s = []
3     for i in range(n): #catching a sample of 30 fish
4         s.append(my_fish())
5     s = np.array(s)
6     std = s.std() #sample standard deviation
7     mean = s.mean() #sample mean
8     a = abs(mean - u_0)
9     prob = 2*(stats.norm(u_0, std/math.sqrt(n)).cdf(u_0 - a))
10    if prob < 0.1:
11        return 1 #rejected
12    else:
13        return 0 #accepted

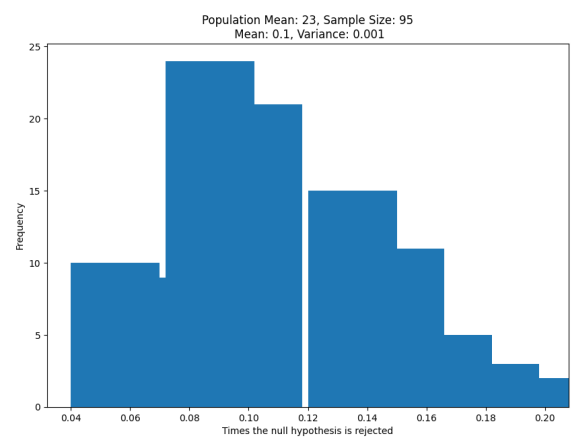
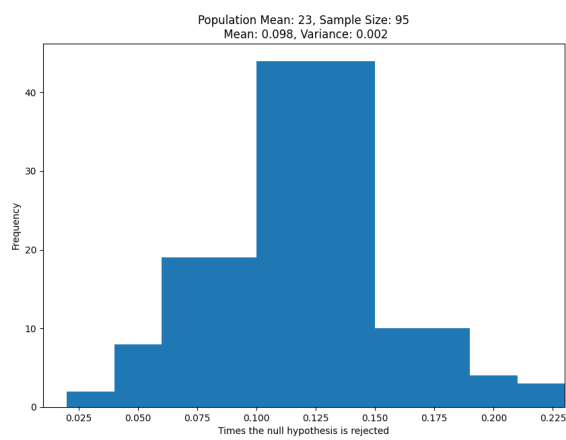
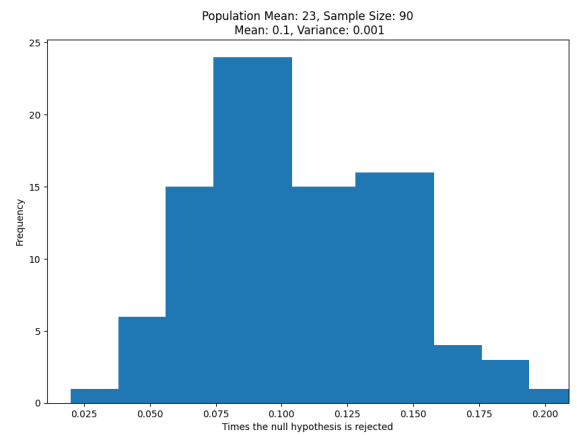
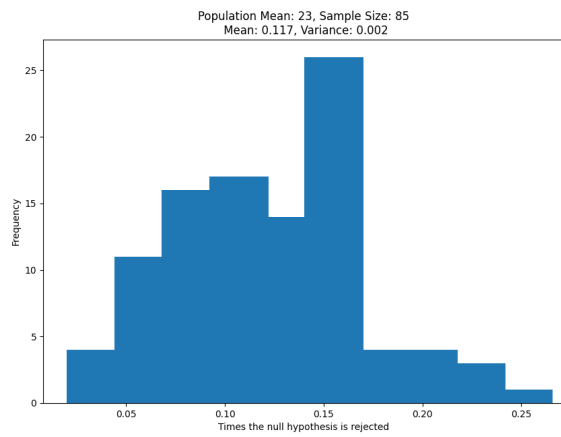
```

We then use the following function which calls the `testing()` function, performs 50 experiments and uses the average of their outcome to calculate a single expected value. It repeats this for 100 iterations and plots a histogram of the expected values.

```

1 def simulation():
2     u_0 = 23 #population mean
3     n = 40
4     expected = []
5     for j in range(100):
6         outcomes = []
7         for i in range(50):
8             outcomes.append(testing(u_0,n))
9         expected.append(sum(outcomes)/50)
10    #plotting histogram
11    fig, ax = plt.subplots(figsize =(10, 7))
12    a = np.arange(0,2,0.1)
13    ax.hist(expected, bins = a)
14    ax.set_xlabel('Times the null hypothesis is rejected')
15    ax.set_ylabel('Frequency')
16    plt.title('Population Mean: '+str(u_0)+' , Sample Size: '+str(n)+'\n Mean: '+str(round(
17        statistics.mean(expected),3))+', Variance: '+str(round(statistics.variance(expected),3)))
18    plt.savefig("Q5_histograms/Q5.2.3_"+str(n)+"(3).png")
19    plt.show()

```



Conducting simulations with various values of  $n$ , we achieve the above histograms. In order to determine the least value of  $n$  to ensure that the null hypothesis is not wrongly rejected more than 10 percent of the times, we want a value of  $n$  that gives a mean of 0.1. According to the simulations, the value of  $n$  satisfying this condition is **90**.