

# Notebook

January 27, 2020



### 0.0.1 Question 1

Let  $x_1, x_2, \dots, x_n$  be a list of numbers. You can think of each index  $i$  as the label of a household, and the entry  $x_i$  as the annual income of Household  $i$ . Define the *mean* or *average* of the list to be  $\mu = \frac{1}{n} \sum_{i=1}^n x_i$ .

**Question 1a)** The  $i$ th *deviation from average* is the difference  $x_i - \mu$ . In Data 8 you saw in numerical examples that the [sum of all these deviations is 0](#). Now prove that fact. That is, show that  $\sum_{i=1}^n (x_i - \mu) = 0$ .

**Note:** In this class, you must always put your answer in the cell that immediately follows the question. DO NOT create any cells between this one and the one that says *Write your answer here, replacing this text*.

$$\sum_{i=1}^n x_i - \sum_{i=1}^n \mu = 0$$

$$\sum_{i=1}^n x_i = \sum_{i=1}^n \mu = 0$$

$$\sum_{i=1}^n x_i = n\mu$$

$$\sum_{i=1}^n x_i = n \frac{1}{n} \sum_{i=1}^n x_i$$

$$\sum_{i=1}^n x_i = \sum_{i=1}^n x_i$$

$$\sum_{i=1}^n x_i - \sum_{i=1}^n x_i = 0$$

$$\sum_{i=1}^n (x_i - \mu) = 0$$



**Question 1b)** Recall that the *variance* of a list is defined as the *mean squared deviation from average*, and that the *standard deviation* (SD) of the list is the square root of the variance. The SD is in the same units as the data and measures the rough size of the deviations from average.

Denote the variance of the list by  $\sigma^2$ . Write a math expression for  $\sigma^2$ . We recommend building your expression by reading the definition of variance from right to left. That is, start by writing the notation for "average", then "deviation from average", and so on.

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$$



**Question 1c)** Suppose you have to predict the value of  $x_i$  for some  $i$ , but you don't get to see  $i$  and you certainly don't get to see  $x_i$ . You decide that whatever  $x_i$  is, you're just going to use your favorite number as your predictor, and your favorite number is  $\mu$ .

The *error* in your prediction is  $x_i - \mu$ , which is your old friend the deviation from average. Thus the *mean squared error* (MSE) of your predictor  $\mu$  over the entire list is the mean squared deviation from average, which is your old friend the variance. So we will write  $\sigma^2 = MSE(\mu)$ .

Now suppose I decide that whatever  $x_i$  is, I'm just going to use *my* favorite number as my predictor, and my favorite number is  $c$ . Write a math expression for  $MSE(c)$ . Again, go from right to left: first  $c$ , then the error, and so on.

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - c)^2}{n}$$





**Question 1d)** Whose predictor is better? It seems reasonable to guess that your predictor  $\mu$  is better than my favorite but possibly weird  $c$ . Show that  $MSE(c) > MSE(\mu)$  for all  $c \neq \mu$ , by the method indicated below.

- Write the error  $x_i - c$  as  $x_i - c = (x_i - \mu) + (\mu - c)$ .
- Substitute this expression for  $x_i - c$  in your formula for  $MSE(c)$ .
- Expand the square and use properties of sums; don't forget what you showed in Part **a**.

This shows that  $\mu$  is the *least squares* constant predictor. In Data 8 you found (numerically) the [least squares linear predictor](#) of a variable  $y$  based on a related variable  $x$ . We will return to that later in this course, using a generalization of the calculation in this exercise.

$$\begin{aligned}
 x_i - c &= (x_i - \mu) + (\mu - c) \\
 MSE(c) &= \frac{1}{n} \sum_{i=1}^n (x_i - c)^2 \\
 &= \frac{1}{n} \sum_{i=1}^n ((x_i - \mu) + (\mu - c))^2 \\
 &= \frac{1}{n} \sum_{i=1}^n ((x_i - \mu)(x_i - \mu) + (x_i - \mu)(\mu - c) + (\mu - c)(x_i - \mu) + (\mu - c)(\mu - c)) \\
 &= \frac{1}{n} \sum_{i=1}^n ((x_i - \mu)^2 + 2(x_i - \mu)(\mu - c) + (\mu - c)^2) \\
 &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 + \frac{1}{n} \sum_{i=1}^n (2(x_i - \mu)(\mu - c)) + \frac{1}{n} \sum_{i=1}^n (\mu - c)^2
 \end{aligned}$$



### 0.0.2 Question 2

Consider the function  $f(x) = x^2$  for  $-\infty < x < \infty$ .

**Question 2a)** Find the equation of the tangent line to  $f$  at  $x = 0$ .

$$f'(x) = 2x$$

$$f'(0) = 2(0)$$

$$f'(x) = 0$$

$$f(0) = (0)^2$$

$$f(0) = 0$$

$$y - y_1 = m(x - x_1)$$

$$y - 0 = 0(x - 0)$$

$$y = 0$$



**Question 2b)** Find the equation of the tangent line to  $f$  at  $x = 8$ .

$$f'(8) = 16$$

$$f(8) = 8^2$$

$$f(8) = 64$$

$$y - y_1 = m(x - x_1)$$

$$y - 64 = 16(x - 8)$$

$$y = 16x - 64$$



**Question 2c)** Write code to plot the function  $f$ , the tangent line at  $x = 8$ , and the tangent line at  $x = 0$ . Set the range of the x-axis to  $(-15, 15)$  and the range of the y-axis to  $(-100, 300)$  and the figure size to  $(4,4)$ .

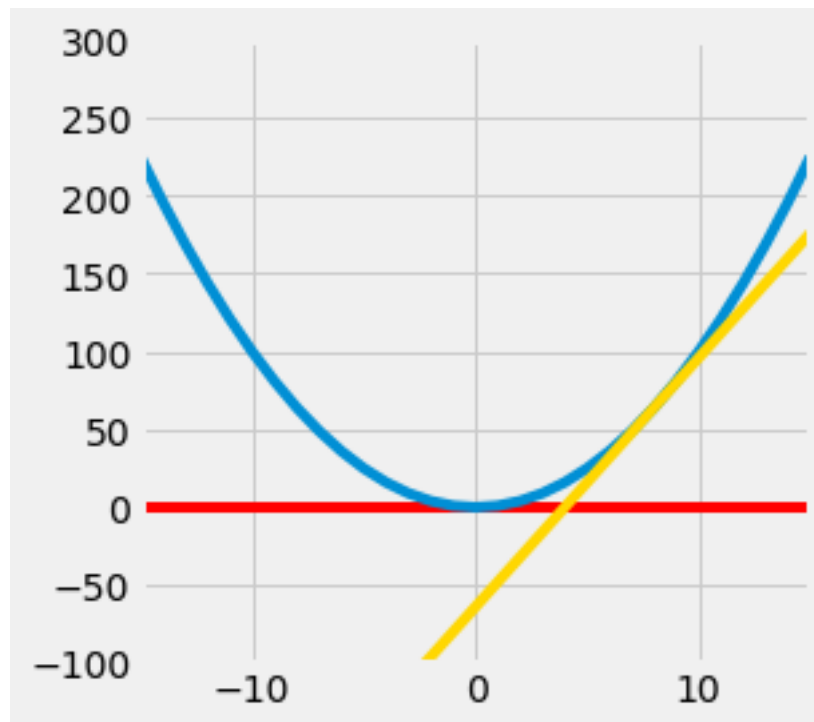
Your resulting plot should look like this:

You should use the `plt.plot` function to plot lines. You may find the following functions useful:

- `plt.plot(..)`
- `plt.figure(figsize=..)`
- `plt.ylim(..)`
- `plt.axhline(..)`

```
In [9]: def f(x):  
        return x**2  
  
        def df(x):  
            return 16*x  
  
        def plot(f, df):  
            x = np.arange(-15, 16)  
            plt.figure(figsize = (4, 4))  
            plt.axhline(0, color = 'red')  
            plt.ylim(bottom = -100, top = 300)  
            plt.xlim(left = -15, right = 15)  
            plt.plot(x, f(x))  
            plt.plot(x, df(x) - 64, color = 'gold')
```

```
plot(f, df)
```





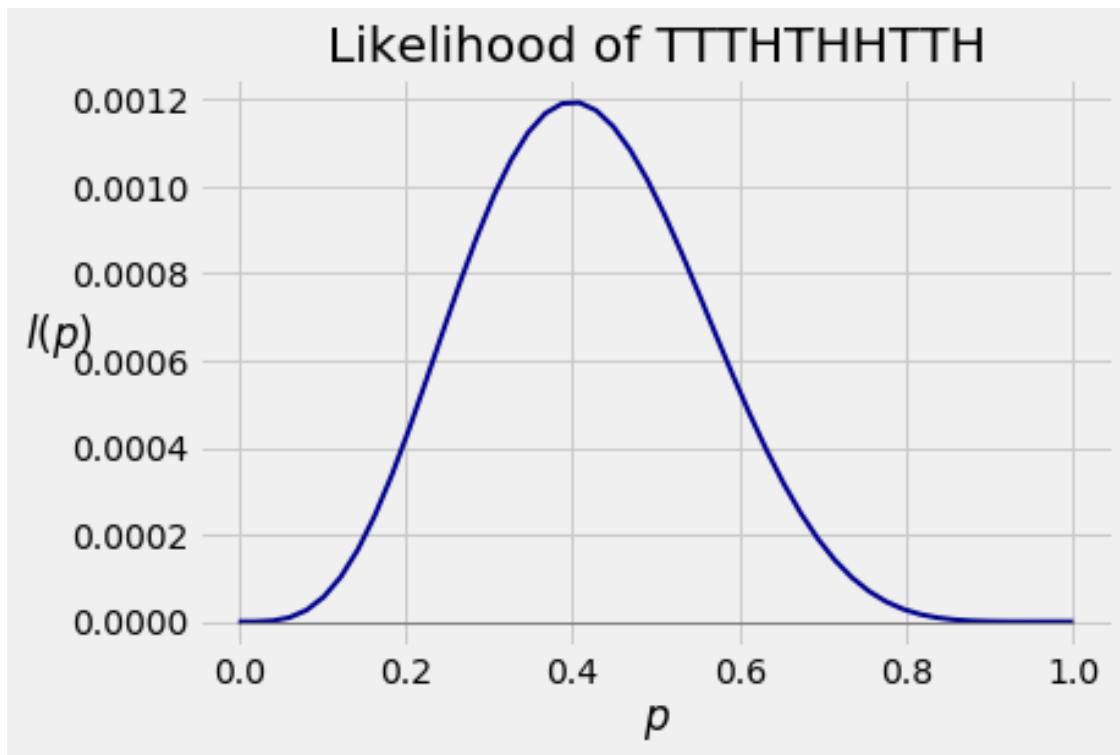


**Question 3b)** I have a coin that lands heads with an unknown probability  $p$ . I toss it 10 times and get the sequence TTTHTHHTTH.

If you toss this coin 10 times, the chance that you get the sequence above is a function of  $p$ . That function is called the *likelihood* of the sequence TTTHTHHTTH, so we will call it  $l$ .

Plot the graph of  $l$  as a function of  $p$  for  $p \in [0, 1]$ .

```
In [12]: p = np.linspace(0, 1)
likelihood = p**4 * (1 - p)**6
plt.plot(p, likelihood, lw=2, color='darkblue') # lw is line width
plt.plot([0, 1], [0, 0], lw=1, color='grey')   # horizontal axis
plt.xlabel('$p$')
plt.ylabel('$l(p)$', rotation=0)
plt.title('Likelihood of TTTHTHHTTH');
```





**Question 3c)** The value  $\hat{p}$  at which the likelihood function attains its maximum is called the *maximum likelihood estimate* (MLE) of  $p$ . Among all values of  $p$ , it is the one that makes the observed data most likely.

Please provide the value of  $\hat{p}$  and also a simple interpretation of that value in terms of the data TTTHTHTTH.

$$\hat{p} = 0.4$$

0.4 is the probability that flipping a coin will result in a head. In terms of the data, having a 0.4 probability of flipping a head will have a 0.0012 probability of giving you the sequence TTTHTHTTH.



**Question 3d)** Explain why the value  $\hat{p}$  at which the function  $l$  attains its maximum is the same as the value at which the function  $\log(l)$  attains its maximum. To clarify,  $\log(l)$  is the composition of  $\log$  and  $l$ :  $\log(l)$  at  $p$  is  $\log(l(p))$ . Even though it doesn't make a difference for this problem,  $\log$  is now and forevermore the log to the base  $e$ , not to the base 10.

It might help to compare  $\log(x_1)$  and  $\log(x_2)$  for  $x_1 < x_2$ .

The observation in this exercise is hugely important in data science because many probabilities are products and the log function turns products into sums. It's much simpler to work with a sum than with a product.

The value  $\hat{p}$  at which the function  $l$  attains its maximum is the same as the value at which  $\log(l)$  attains its maximum is because the y-values on the graph of function  $l$  at each  $x$  are essentially data points. By applying the log function to each data point, the maximum at  $\hat{p} = 0.4$  will be the maximum for  $\log(l)$  since log functions increase infinitely in the positive direction, so the largest  $l(p)$  will result in the maximum for  $\log(l)$ .



**Question 3e)** Use Part c and calculus to find  $\hat{p}$ . Using Part d makes the calculus much easier. You don't have to check that the value you've found produces a max and not a min – we'll spare you that step.

$$l(p) = p^4 * (1 - p)^6$$

$$\log(l(p)) = \log(p^4 * (1 - p)^6)$$

$$\log(l(p)) = \log(p^4) + \log((1 - p)^6)$$

$$d/dx \log(l(p)) = 4p^3/(p^4) - 6(1 - p)^5/(1 - p) * 6$$

$$d/dx \log(l(p)) = 4/p - 6/(1 - p)$$

$$0 = 4/p - 6/(1 - p)$$

$$6/(1 - p) = 4/p$$

$$6p = 4 - 4p$$

$$\hat{p} = 0.4$$





### 0.0.3 Question 4

Data science is a rapidly expanding field and no degree program can hope to teach you everything that will be helpful to you as a data scientist. So it's important that you become familiar with looking up documentation and learning how to read it.

Below is a section of code that plots a three-dimensional "wireframe" plot. You'll see what that means when you draw it. Replace each `# Your answer here` with a description of what the line above does, what the arguments being passed in are, and how the arguments are used in the function. For example,

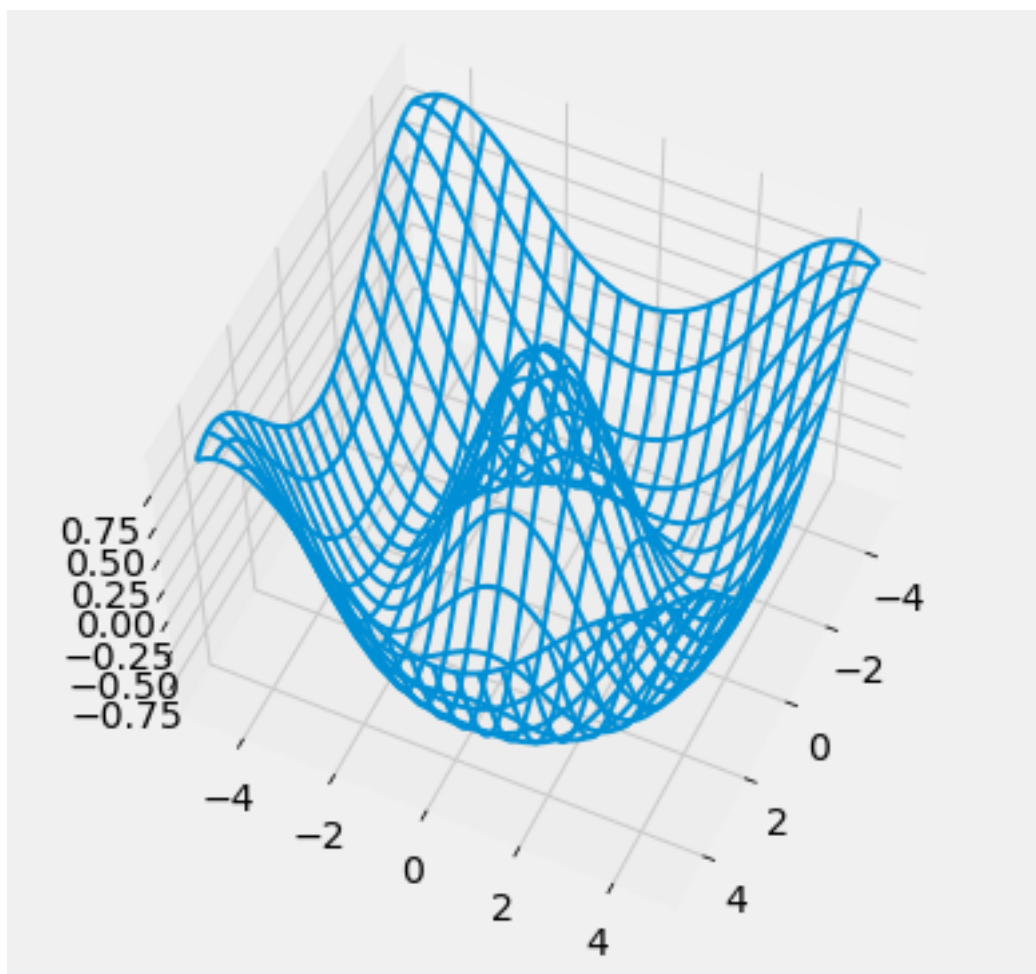
```
np.arange(2, 5, 0.2)
# This returns an array of numbers from 2 to 5 with an interval size of 0.2
```

**Hint:** The Shift + Tab tip from earlier in the notebook may help here. Remember that objects must be defined in order for the documentation shortcut to work; for example, all of the documentation will show for method calls from `np` since we've already executed `import numpy as np`. However, since `z` is not yet defined in the kernel, `z.reshape()` will not show documentation until you run the line `z = np.cos(squared)`.

```
In [13]: from mpl_toolkits.mplot3d import axes3d

u = np.linspace(1.5*np.pi, -1.5*np.pi, 100)
# This returns 100 evenly spaced numbers over the interval [-1.5*np.pi, 1.5*np.pi]
[x,y] = np.meshgrid(u, u)
# This takes in 1D arrays and uses them as coordinates to return two 2D arrays
squared = np.sqrt(x.flatten()**2 + y.flatten()**2)
z = np.cos(squared)
# This returns the cos function applied to a 1D array of the flattened array 'squared'
z = z.reshape(x.shape)
# This returns the data in the z array but in the new shape of array x

fig = plt.figure(figsize=(6, 6))
ax = fig.add_subplot(111, projection='3d')
# This adds a subplot to the current figure with 1 row, 1 column, and index starting at 1 with
ax.plot_wireframe(x, y, z, rstride=5, cstride=5, lw=2)
# This plots a 3D wireframe with x, y, z data points with row stride of 5, column stride of 5,
ax.view_init(elev=60, azim=25)
# This sets the elevation angle of ax to 60 and the azimuth angle to 25
plt.savefig("figure1.png")
# This saves the current figure as "figure1.png"
```



#### 0.0.4 Question 5

Much of data analysis involves interpreting proportions – lots and lots of related proportions. So let's recall the basics. It might help to start by reviewing [the main rules](#) from Data 8, with particular attention to what's being multiplied in the multiplication rule.

**Question 5a)** The Pew Research Foundation publishes the results of numerous surveys, one of which is about the [trust that Americans have](#) in groups such as the military, scientists, and elected officials to act in the public interest. A table in the article summarizes the results.

Pick one of the options (i) and (ii) to answer the question below; if you pick (i), fill in the blank with the percent. Then, explain your choice.

The percent of surveyed U.S. adults who had a great deal of confidence in both scientists and religious leaders

(i) is equal to \_\_\_\_\_.

(ii) cannot be found with the information in the article.

(ii) because we do not know if the Americans who votes in the survey were able to vote for multiple choices, so Americans who voted for scientists may have, or may not have, voted for religious leaders as well. This means we cannot know for sure if the events are mutually exclusive.



**Question 5d)** (This part is a continuation of the previous two.) Pick all of the options (i)-(iv) that are true for all values of  $p$ . Explain by algebraic or probabilistic reasoning; you are welcome to use your function `no_disease_given_negative` to try a few cases numerically. Your explanation should include the reasons why you *didn't* choose some options.

$P(N \mid T_N)$  is

- (i) equal to 0.95.
- (ii) equal to  $0.999 \times 0.95$ .
- (iii) greater than  $0.999 \times 0.95$ .
- (iv) greater than 0.95.

(iii, iv) This is because  $P(N \mid T_N)$  given a value  $p$  in the range 0 to 1 is  $(0.999 * 0.95) / ((0.999 * 0.95) + (0.001 * (1-p)))$ . Looking at this equation, we see that whatever  $p$  is, the value multiplied by 0.001 will essentially result in an answer close to or equal to 1, so it will be greater than choices (iii) and (iv). Choices (i) and (ii) are not true because  $p$  could be any number chosen in the range of (0, 1), so it cannot be a constant value.



**Question 5e)** Suzuki is one of most commonly owned makes of cars in our county (Alameda). A car heading from Berkeley to San Francisco is pulled over on the freeway for speeding. Suppose I tell you that the car is either a Suzuki or a Lamborghini, and you have to guess which of the two is more likely.

What would you guess, and why? Make some reasonable assumptions and explain them (data scientists often have to do this), justify your answer, and say how it's connected to the previous parts.

I would guess that the Lamborghini would be the car that it pulled over for speeding. Drawing a similar tree diagram to the disease diagram, disease would be replaced with speeding, and no disease would be replaced with not speeding. Test positive would be Lamborghini, and test negative would be Suzuki. Using similar probabilities, of the cars caught speeding, a high end luxury sports car would have a much higher probability of being caught speeding over a lower end car. This also makes sense logically as Lamborghinis and sport cars catch the eyes of police officers much more than an average Suzuki would.



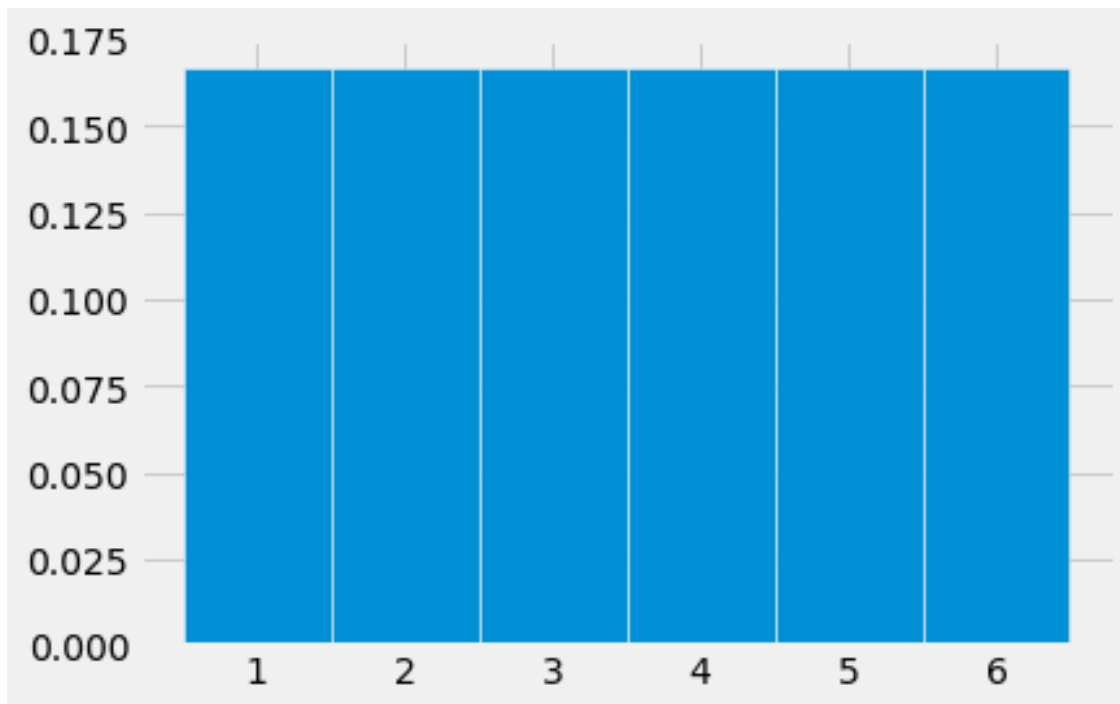


**Question 6a)** Define a function `integer_distribution` that takes an array of integers and draws the histogram of the distribution using unit bins centered at the integers and white edges for the bars. The histogram should be drawn to the density scale. The left-most bar should be centered at the smallest integer in the array, and the right-most bar at the largest.

Your function does not have to check that the input is an array consisting only of integers. The display does not need to include the printed proportions and bins.

If you have trouble defining the function, go back and carefully read all the lines of code that resulted in the probability histogram of the number of spots on one roll of a die. Pay special attention to the bins.

```
In [23]: def integer_distribution(x):  
        min_value = min(x)  
        max_value = max(x)  
        unit_bins = np.arange(min_value - 0.5, max_value + 0.6)  
        plt.hist(x, bins = unit_bins, ec='white', density=True)  
  
        integer_distribution(faces)
```

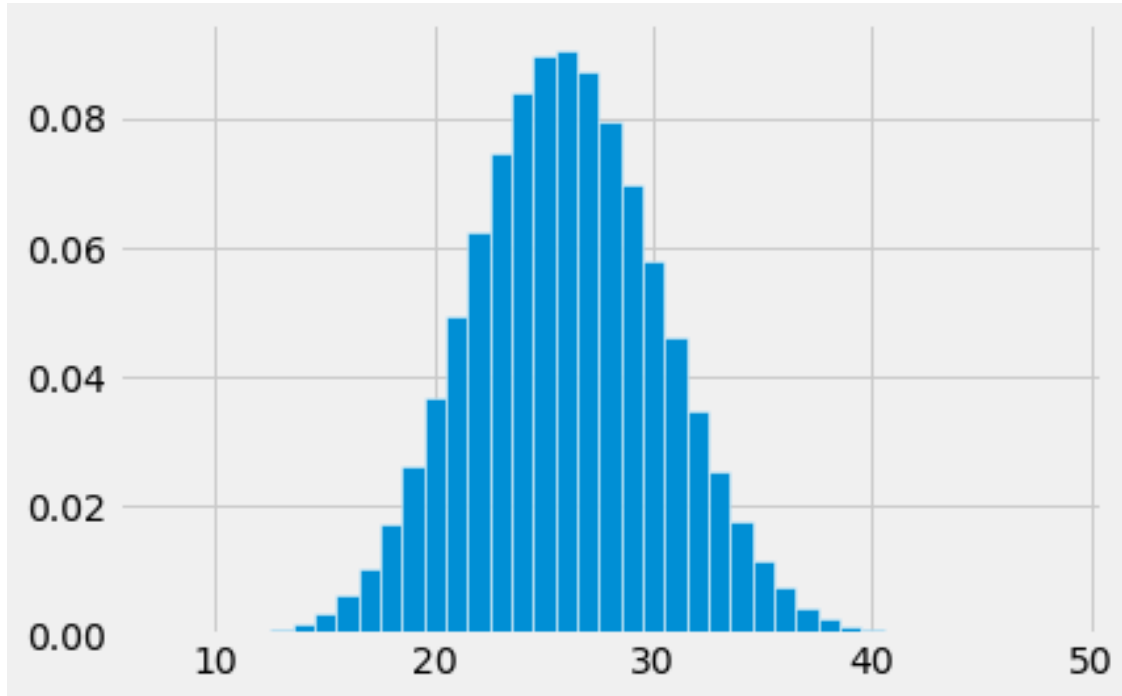




**Question 6c)** Replace the “...” in the code cell below with a Python expression so that the output of the cell is an empirical histogram of 500,000 simulated counts of black people in 100 draws made at random with replacement from the population eligible for Swain’s jury panel.

After you have drawn the histogram, you might want to take a moment to recall the conclusion reached in Data 8 based on the data that Swain’s panel had 8 black people in it.

```
In [26]: simulated_counts = np.random.multinomial(100, [.74, .26], 500000)[: , -1]
         integer_distribution(simulated_counts)
```





**Question 6d)** As you know, the count of black people in a sample of 100 people drawn at random from the eligible population is expected to be 26. Just by looking at the histogram in Part **c**, and **no other calculation**, pick the correct option and explain your choice. You might want to refer to the [Data 8 textbook](#) again.

The SD of the distribution of the number of black people in a random sample of 100 people drawn from the eligible population is closest to

(i) 1.4

(ii) 4.4

(iii) 7.4

(iv) 10.4

(ii) because the bell curve starts to become a "right-way-up cup" around 30 and 31, so 4.4 is the closest to being the SD of the distribution of the number of black people in a random sample.



**Question 6e)** The *normal curve with mean  $\mu$  and SD  $\sigma$*  is defined by

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}, \quad -\infty < x < \infty$$

Redraw your histogram from Part **c** and overlay the normal curve with  $\mu = 26$  and  $\sigma$  equal to the choice you made in Part **d**. You just have to call `plt.plot` after `integer_distribution`. Use `np.e` for `e`. For the curve, use 2 as the line width, and any color that is easy to see over the blue histogram. It's fine to just let Python use its default color.

Now you can see why centering the histogram bars over the integers was a good idea. The normal curve peaks at 26, which is the center of the corresponding bar.

```
In [27]: mu = 26
         sigma = 4.4
         x = np.linspace(0, 50, 200)
         f_x = 1/(np.sqrt(2*np.pi)*sigma) * np.e**((-1/2)*((x-mu)/sigma)**2)
         integer_distribution(simulated_counts)
         plt.plot(x, f_x)
```

```
Out[27]: [<matplotlib.lines.Line2D at 0x7fc2441636d8>]
```

