

Part D) Are X and Y independent or dependent? Fully justify your answer in the cell below using LaTeX and the mathematical definition of independence.

We know that if $P(X = x, Y = y) = P(X = x) \cdot P(Y = y)$, then X and Y are independent.

First we figure out $P(X = 1, Y = 1) = P(X = 1) \cdot P(Y = 1)$

$$\frac{1}{3} \neq \frac{5}{12} \cdot \frac{7}{12}$$

$$P(X = 1, Y = 2) = P(X = 1) \cdot P(Y = 2)$$

$$\frac{1}{12} \neq \frac{5}{12} \cdot \frac{5}{12}$$

$$P(X = 2, Y = 1) = P(X = 2) \cdot P(Y = 1)$$

$$\frac{1}{6} \neq \frac{1}{6} \cdot \frac{7}{12}$$

$$P(X = 2, Y = 2) = P(X = 2) \cdot P(Y = 2)$$

$$0 \neq \frac{1}{6} \cdot \frac{5}{12}$$

$$P(X = 4, Y = 1) = P(X = 4) \cdot P(Y = 1)$$

$$\frac{1}{12} \neq \frac{5}{12} \cdot \frac{7}{12}$$

$$P(X = 4, Y = 2) = P(X = 4) \cdot P(Y = 2)$$

$$\frac{1}{12} \neq \frac{5}{12} \cdot \frac{5}{12}$$

Thus, we know that all these possibilities do not equal each other, thus we know that they are dependent of each other.

Part A) If $\text{Cov}(X, Y) = 0$, what does this tell us about the random variables X and Y ?

If $\text{Cov}(X, Y) = 0$, it tells us that X and Y are uncorrelated, and there is no linear relationship between the two variables.

Part B) Given the following joint pmf for discrete random variables X and Y :

	$Y = 0$	$Y = 1$	$Y = 2$
$X = 0$	$\frac{1}{6}$	$\frac{1}{4}$	$\frac{1}{8}$
$X = 1$	$\frac{1}{8}$	$\frac{1}{6}$	$\frac{1}{6}$

- i). Calculate $Cov(X, Y)$.
- ii). Calculate $\rho(X, Y)$

Show all steps for both parts using Markdown and LaTeX in the cell below:

i). Calculate $Cov(X, Y)$.

$$E[X] = 0 \cdot \left(\frac{1}{6} + \frac{1}{4} + \frac{1}{8}\right) + 1 \cdot \left(\frac{1}{8} + \frac{1}{6} + \frac{1}{6}\right) = \frac{11}{24}$$

$$E[Y] = 0 \cdot \left(\frac{1}{6} + \frac{1}{8}\right) + 1 \cdot \left(\frac{1}{4} + \frac{1}{6}\right) + 2 \cdot \left(\frac{1}{8} + \frac{1}{6}\right) = 1$$

$$E[XY] = (0 \cdot 0)\left(\frac{1}{6}\right) + (1 \cdot 0)\left(\frac{1}{8}\right) + (0 \cdot 1)\left(\frac{1}{4}\right) + (1 \cdot 1)\left(\frac{1}{6}\right) + (0 \cdot 2)\left(\frac{1}{8}\right) + (1 \cdot 2)\left(\frac{1}{6}\right) = \frac{1}{2}$$

$$Cov(X, Y) = E[XY] - E[X]E[Y] = \frac{1}{2} - \frac{11}{24} \cdot 1 = \frac{1}{24}$$

ii). Calculate $\rho(X, Y)$

$$\rho(X, Y) = \frac{Cov(X, Y)}{SD(X)SD(Y)}$$

$$SD(X) = \sqrt{Var(X)} = \sqrt{E(X^2) - (E(X))^2} = \sqrt{(0^2 \cdot (\frac{1}{6} + \frac{1}{4} + \frac{1}{8}) + 1^2 \cdot (\frac{1}{8} + \frac{1}{6} + \frac{1}{6})) - (\frac{11}{24})^2} = \sqrt{\frac{143}{576}}$$

$$SD(Y) = \sqrt{Var(Y)} = \sqrt{E(Y^2) - (E(Y))^2} = \sqrt{(0^2 \cdot (\frac{1}{6} + \frac{1}{8}) + 1^2 \cdot (\frac{1}{4} + \frac{1}{6}) + 2^2 \cdot (\frac{1}{8} + \frac{1}{6})) - (1)^2} = \sqrt{\frac{7}{12}}$$

$$\rho(X, Y) = \frac{\frac{1}{24}}{\sqrt{\frac{143}{576}} \cdot \sqrt{\frac{7}{12}}} = \frac{2 \cdot \sqrt{3003}}{1001} \approx 0.10949$$

Part C) This part is **NOT** related to the parts above.

Suppose you're only given the following information about two joint random variables X and Y :

$$\mu_X = 6, \quad \mu_Y = 5, \quad \sigma_X^2 = 4, \quad \sigma_Y^2 = 9 \text{ and } E[XY] = 27$$

.

For each of the quantities below, calculate if you have enough information, showing all steps. If not, explain what additional info you'd need.

i). $Cov(X, Y)$

ii). $Cov(Y, X)$

iii). $\rho(X, Y)$

Answer all parts in the ONE markdown cell below, fully justifying your answer:

i). $Cov(X, Y)$

$$Cov(X, Y) = E[XY] - \mu_X \cdot \mu_Y = 27 - 6 \cdot 5 = -3$$

ii). $Cov(Y, X)$

$$\text{If } Cov(Y, X) = E[YX] - \mu_Y \cdot \mu_X = E[XY] - \mu_X \cdot \mu_Y = Cov(X, Y) = 27 - 6 \cdot 5 = -3$$

iii). $\rho(X, Y)$

$$\rho(X, Y) = \frac{Cov(X, Y)}{SD(X)SD(Y)} = \frac{Cov(X, Y)}{\sqrt{\sigma_X^2} \cdot \sqrt{\sigma_Y^2}} = \frac{-3}{2 \cdot 3} = \frac{-1}{2}$$

Back to top

0.1 (2 pts) Problem 4

If we're trying to predict the results of the Clinton vs. Trump 2016 presidential race:

- i). What is the population of interest?
- ii). What is the sampling frame?

Give both of your answers in the same below in Markdown.

- i) The population of interest is the eligible voting population in the U.S. in the 2016 presidential election.
- ii) The sampling frame is all the list of the registered voters in the U.S.

Back to top

0.2 Problem 5 (11 pts)

Part A For your convenience, the actual results of the vote in the four pivotal states is repeated below:

State	% Trump	% Clinton	Total Voters
florida	49.02	47.82	9,419,886
michigan	47.50	47.27	4,799,284
pennsylvania	48.18	47.46	6,165,478
wisconsin	47.22	46.45	2,976,150

Using the table above, write a function `draw_state_sample(N, state)` that returns a sample with replacement of N voters from the given state, using the percentages given in the table above. Your result should be returned as a list, where the first element is the number of Trump votes, the second element is the number of Clinton votes, and the third is the number of Other votes. For example, `draw_state_sample(1500, "florida")` could return `[727, 692, 81]`. You may assume that the state name is given in all lower case.

Hint: You might find `np.random.multinomial` useful.

```
In [29]: vote = {
    "florida" : [0.4902,0.4782,0.0316],
    "michigan": [0.4750, 0.4727, 0.0523],
    "pennsylvania": [0.4818, 0.4746, 0.0436],
    "wisconsin": [0.4722, 0.4645, 0.0633],
    }

    def draw_state_sample(N, state):
        percent = vote[state]
        sample = np.random.multinomial(N,percent)
        return sample

    s = draw_state_sample(1500,"florida")
    print(s)
```

```
[754 695  51]
```

```
In [30]: grader.check("q5a")
```

```
Out[30]: q5a results: All test cases passed!
```

Part D i). Make a **frequency** histogram of **simulations**. This is a histogram of the sampling distribution of Trump's proportion advantage in Pennsylvania.

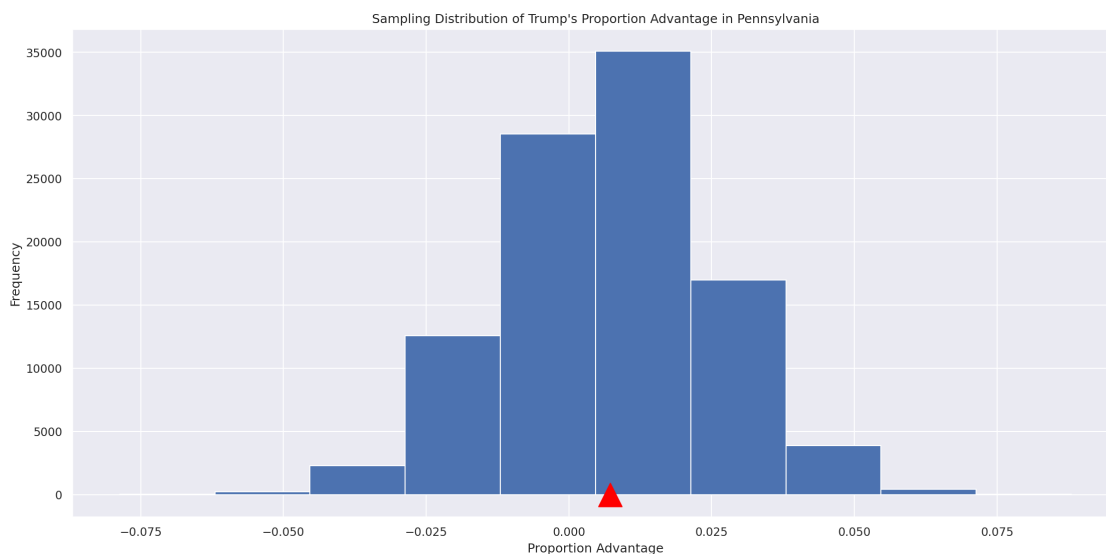
Hint: You should use the `plt.hist` function in your code.

Make sure to include a title as well as axis labels. You can do this using `plt.title`, `plt.xlabel`, and `plt.ylabel`.

ii). Based on your simulation, what is the probability that a random sample of 1500 will correctly predict that Trump wins Pennsylvania? (i.e. what proportion of these simulations predict a Trump victory?) Assign your answer to `prob_penn_1500_random_correct`

```
In [35]: # Part (i):
plt.hist(simulations)
plt.title("Sampling Distribution of Trump's Proportion Advantage in Pennsylvania")
plt.xlabel("Proportion Advantage")
plt.ylabel("Frequency")
mean = np.mean(simulations)
# your code for the histogram above here. The code below plots a red marker at the mean:
plt.scatter(mean, -1, marker='^', color='red', s=500)
plt.show
```

```
Out[35]: <function matplotlib.pyplot.show(close=None, block=None)>
```



```
In [36]: # Part (ii):
```

```
correct = sum(win > 0 for win in simulations)
prob_penn_1500_random_correct = correct/len(simulations)

prob_penn_1500_random_correct
```

Out[36]: 0.64716

Back to top

0.3 Problem 6 (10 pts)

Throughout this problem, adjust the selection of voters so that there is a 0.5% bias in favor of Clinton in each of these states.

For example, in Pennsylvania, Clinton received 47.46% of the votes and Trump 48.18%. Increase the population of Clinton voters to $47.46\% + 0.5\%$ and correspondingly decrease the percent of Trump voters.

Part A Simulate Trump's advantage across 100,000 simple random samples of 1500 voters for the **state of Pennsylvania** and store the results of each simulation in an `np.array` called `biased_simulations`.

That is, `biased_simulation[i]` should hold the result of the `i+1`th simulation.

That is, your answer to this problem should be just like your answer from Question 5C, but now using samples that are biased as described above.

```
In [42]: def draw_biased_state_sample(N, state):
        vote_bias = {
            'florida': [0.4852, 0.4832, 0.0316],
            'michigan': [0.47, 0.4777, 0.0523],
            'pennsylvania': [0.4768, 0.4796, 0.0436],
            'wisconsin': [0.4672, 0.4675, 0.0633]
        }
        return np.random.multinomial(N, vote_bias[state])
```

```
biased_simulations = [trump_advantage(draw_biased_state_sample(1500, "pennsylvania")) for i in
```

```
In [43]: grader.check("q6a")
```

```
Out[43]: q6a results: All test cases passed!
```

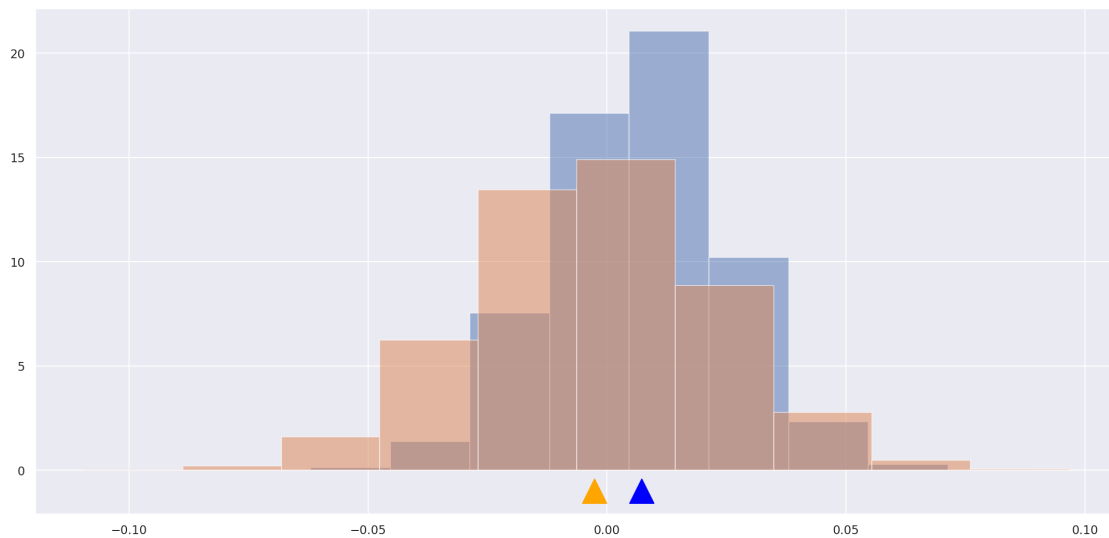

Part B Create a plot of **overlaid DENSITY** histograms of the following: - The new sampling distribution of Trump's proportion advantage in Pennsylvania using these biased samples - The sampling distribution of the unbiased samples from Problem 5D (plotted as a density, not a frequency histogram)

Include 2 markers (of different colors) with the sample means for each distribution (see 5D for code how to do this). The colors of the markers should correspond to the colors of the density histograms.

Make sure to give your plot a title, label the x and y axes and include a legend. Use the parameter **alpha** to adjust the transparency of each histogram.

```
In [44]: plt.hist(simulations, density=True, alpha=0.5)
plt.hist(biased_simulations, density=True, alpha=0.5)
mean_s = np.mean(simulations)
biased = np.mean(biased_simulations)
plt.scatter(mean_s, -1, marker='^', color='blue', s=500)
plt.scatter(biased, -1, marker='^', color='orange', s=500)
```

Out[44]: <matplotlib.collections.PathCollection at 0x7f55da864760>



Summarize the findings from these simulations:

i). Based on your simulations, what was the **chance of error** in correctly predicting that Trump wins using the **unbiased** samples of 1500 people from each state? Many people, even well educated ones, assume that this number should be 0%. After all, how could a non-biased sample be wrong? Give a mathematical explanation as to why it isn't 0% (or close to 0%). This is the type of incredibly important intuition we hope to develop in you throughout this class and your future data science coursework.

ii). What was the chance of error in predicting the results using the **biased** samples and how different is it from your answer in part(i)? Recall, we only biased the samples by 0.5%. However, even a bias this small in the percentages can lead to a much larger chance of error in prediction of the final result.

i) The chance of error is approximately 30%. Its because there is variability in the random samples. Even if we have an unbiased sample, there will be random variation in sample results. Thus, due to the randomness of sampling, there may be a possibility that a sample might have distribution of voters that doesn't match the population proportion. .

ii) The chance of error in predicting results using biased samples were 50%. This is because when the samples were biased by 0.5%, the chance of error in predicting the result increased. This is because the bias shifted the sample results away from the true population proportions. Since the bias is towards Clinton, the chance of error for Trump is higher. Thus even a bias this small can lead to significant impact to error in prediction.

Part B Compare your observations from 7a to your observations in 6d. Did the chance of error increase or decrease in each case and why? What do these changes imply about the impact of sample size on the sampling error and on the bias?

It seems by increasing the sample size, the unbiased sample's incorrect proportions were reduced, which may imply that larger sample sizes lead to accurate predictions. But the biased basically stayed the same, but did increase by 0.6. Thus, it does show that larger sample sizes may provide more precision and accuracy , but don't really reduce the impact of the bias.

Part C Is it possible to correctly predict Trump's victory with less than 1% error using **unbiased sampling**? Rerun the simulation (in each of the 4 states) with increasing sample sizes and 100,000 simulations to determine if you can find an approximate minimum sample size (it doesn't have to be exact) such that the probability of correctly predicting Trump's victory is at least 99% (assuming your sample is unbiased).

```
In [51]: target = 0.99
        sample = 1500
        unbiased_correct = 0

        while unbiased_correct < target:
            unbiased_correct = np.array([trump_wins(sample) for i in range(100000)]).mean()
            sample += 1000
            if unbiased_correct >= target:
                break
        print("Min Sample size for 99% prediction:", sample, "\nPercentage: " ,unbiased_correct)
        # your code above this line.
        # output the number of samples you used to get to at least 99% accuracy.
```

```
Min Sample size for 99% prediction: 32500
Percentage: 0.9903
```


Part D Is it possible to correctly predict Trump's victory with less than 1% error using **biased sampling**? Use the code cell below to rerun the simulation (in each of the 4 states) with increasing sample sizes. What happens to the probability of error? Explain in the markdown cell below.

It seems that the probability of the error decreases as sample size increases, and converges to some target value, but it cannot achieve less than 1% error due to bias.

```
In [52]: b_target = 0.01
        b_sample = 1500
        biased_correct = 1
        while biased_correct > b_target:
            biased_correct = 1 - np.array([trump_wins_biased(b_sample) for i in range(100000)]).mean()
            b_sample += 1000
            if biased_correct <= target:
                break
        print("Min Sample size for 1% prediction:", b_sample, "\n Percentage: ", biased_correct)
```

```
Min Sample size for 1% prediction: 2500
Percentage: 0.52439
```

