

INSTRUCTIONS

- You have 3 hours to complete the exam.
- The exam is closed book, closed notes, closed computer, and closed calculator, apart from the official final exam reference guide provided with the exam.
- Write each answer **in the space provided for that answer**.
- You can leave all numerical calculations unsimplified.
- Assume that the code statements `import numpy as np` and `from datascience import *` have been executed.

Question 0 (1 point) Write your name and SID in the space provided on one side of every page of the exam.

| | |
|---|--|
| Last name | |
| First name | |
| Student ID number | |
| CalCentral email (<code>_@berkeley.edu</code>) | |
| GSI | |
| Name of the person to your left | |
| Name of the person to your right | |
| <i>All the work on this exam is my own.</i> (please sign) | |

1. (7 points) Income Intervals

Researchers studying annual incomes in a city take a random sample of 400 households and create 10,000 bootstrap samples from the original sample. They then use the bootstrap percentile method to construct an approximate 95% confidence interval for the median household income in the city. This is the method we have always used in Data 8, and you can assume that it works fine in this situation. The 95% confidence interval goes from \$60,000 to \$62,000.

(a) (4 pt) Fill in the bubble next to each statement that must be true based on the information above. **More than one might be true.** No explanations are needed.

- ☐ About 50% of the households in the city have incomes between \$60,000 and \$62,000.
- ☐ About 95% of the households in the city have incomes between \$60,000 and \$62,000.
- ☐ The researchers are estimating that the median household income in the city is between \$60,000 and \$62,000, but they could be wrong.
- ☐ If the researchers had constructed an approximate 90% confidence interval based on the same bootstrap samples they used for the 95% interval, then both ends of their 90% confidence interval would have been inside the range \$60,000 to \$62,000.
- ☐ None of the above is true.

Option 3,4

Common mistakes:

- Option 4 is correct. A 90 percent confidence interval (since you're less confident) is a smaller interval that is contained in your 95 percent interval.
- Option 2 is incorrect because we are only measuring the median household income; not the income of all households in the population

(b) (3 pt) The array `resampled_medians` contains the 10,000 bootstrapped medians. Complete the code below so that the last line evaluates to the left and right endpoints of an approximate 90% confidence interval for the median household income in the city.

```
left_end = _____

right_end = _____

make_array(left_end, right_end)
left_end = percentile(5, resampled_medians)

right_end = percentile(95, resampled_medians)

make_array(left_end, right_end)
```

2. (5 points) Sample Size

The table below shows some percentages under the normal curve, in addition to those already in the exam reference guide.

| Percent in Range | Normal Distribution |
|------------------------|---------------------|
| average \pm 1.3 SDs | about 80% |
| average \pm 1.65 SDs | about 90% |

A researcher wants to estimate the percent of undecided voters in a population, by constructing a confidence interval based on a random sample of voters. Approximately what is the smallest sample size that will ensure

that the width of a 90% confidence interval will be no more than 5%?

90% confidence interval is 1.65 Standard deviations both ways. Width should be no more than 5%, which means that 1.65 standard deviations each way should equate to 5%. $2 * 1.65 * SD(0/1) / \sqrt{sampsiz e} \geq .05$. Note that the standard deviation of a 0,1 population is at most .5, so we can work with that SD for now and notice that any SD smaller will just be a smaller width on top. So, rearranging out terms, we get that the sample size is equal to $(2 * 1.65 * .5 / .05)^2$.

Common mistakes:

- The SDs given in the table are **not** the population SD. These are how much of the interval is contained in a normal distribution within these amounts of SDs away
- The population SD of a 0/1 population (as we are working with) is bounded by a maximum of 0.5
- The formula for the width of a 95% confidence interval appears on your cheat sheet. Notice that we use 4 in the numerator to signify that 2 SDs either way from the average encompasses 95% of the data. In this case, for a 90% confidence interval, we want to look 1.65 standard deviations each way from the average.

3. (10 points) Assessing BOKS

In February of this year, the American Journal of Preventive Medicine presented an analysis of a program called BOKS (Build Our Kids' Success) in Massachusetts. BOKS is a before-school exercise program in which children arrive at school one hour early to engage in games and other physical exercise before the day's classes begin.

Participation in BOKS is voluntary. Researchers studied several hundred BOKS participants in kindergarten through eighth grade, and compared their outcomes with those of children who did not join the BOKS program.

Based on measurements taken after 12 weeks of the BOKS program, a smaller proportion of BOKS participants qualified as obese compared to non-participants. Also, BOKS participants reported feeling deeper social connections to their friends and school and greater happiness and satisfaction with life than they did before the program. Non-participants had no changes to their feelings of well being.

The New York Times reported, "The upshot is that a one-hour, before-school exercise program does seem likely to improve young people's health and happiness, says Dr. Elsie Taveras, a professor at Harvard and head of general pediatrics at Massachusetts General Hospital ..."

(a) (2 pt) Was this an observational study or a randomized controlled experiment? Explain your answer.

This was an observational study, as participation in BOKS was optional. Hence, we are just observing between those who decided to join, and those who did not.

(b) (2 pt) Was there a treatment group? If so, who was in the treatment group?

Yes, there was a treatment group. The kids who participated in BOKS.

(c) (2 pt) Was there a control group? If so, who was in the control group?

Yes, there was a control group. The kids who did not participate in BOKS.

(d) (4 pt) Do you think that the researchers' analysis clearly established the effectiveness of BOKS? If you do, then explain why. If you don't, then give the main reason why not, and provide an alternative explanation for the researchers' results.

No, it did not establish the effectiveness. Since this is an observational study, there might have been confounding factors. An example is noting that kids who joined BOKS might already be living healthier lifestyle, which made them join BOKS in the first place. Hence, they could be more active and eating well outside of the program which caused the outcome we see.

Common Mistakes:

- Just specifying that there were confounding factors was not enough. The question asks for an alternative explanation for the researcher's results, so we are asking for a specific confounding factor
- Many confounding factors exist, though to get credit, it must be relatively specific. It needs to be able to fully explain how the researcher's results came out as such

4. (6 points) Mean Absolute Error

The table `incomes` consists of just one column, labeled `'Income'`. The column contains the incomes of a random sample of workers drawn from a large population. I would like to use this sample to predict the income of a new worker drawn from the population.

To measure how good my prediction is, I will use the mean absolute error instead of the more commonly used mean squared error. For example, if the incomes were \$1, \$2, \$3, and \$4, and I used \$2.50 as my prediction, then the mean absolute error would be

$$\frac{1}{4}(|1 - 2.5| + |2 - 2.5| + |3 - 2.5| + |4 - 2.5|)$$

Help me find the best predictor of income based on the table `incomes`, as follows.

Complete the code below so that the last line evaluates to an array whose first element is the predictor that minimizes the mean absolute error and whose second element is the value of the minimum mean absolute error. You should start by defining a function `mae` that computes the mean absolute error. You are not required to use every blank line.

```
def mae(c):  
  
    return np.average(np.abs(incomes.column(0) - c))  
  
best_predictor = minimize(mae)  
  
mae_of_best_predictor = mae(best_predictor)  
  
make_array(best_predictor, mae_of_best_predictor)
```

Common Mistakes:

- The crux of this problem is noticing that, in order to minimize `mae`, `mae` must take in only one number. It then must work with the `incomes` table inside of the function.
- Passing in a table into the function won't let you minimize it. `Minimize` returns the argument to the function passed in which allows the function to return the smallest value

5. (11 points) Characteristics of Adults

The table `adults` consists of one row for each adult in a population. Each adult has several attributes including age, hours worked per week, employment status, income bracket, and more. The first few rows of `adults` are shown below. The string '`<=50K`' stands for the income bracket “at most \$50,000”.

| Age | Education | Relationship Status | Job Sector | Hours per Week | Country | Income |
|-----|-----------|-----------------------|-------------------|----------------|---------------|--------|
| 39 | Bachelors | Never-married | Adm-clerical | 40 | United-States | <=50K |
| 50 | Bachelors | Married-civ-spouse | Exec-managerial | 13 | United-States | <=50K |
| 38 | HS-grad | Divorced | Handlers-cleaners | 40 | United-States | <=50K |
| 53 | 11th | Married-civ-spouse | Handlers-cleaners | 40 | United-States | <=50K |
| 28 | Bachelors | Married-civ-spouse | Prof-specialty | 40 | Cuba | <=50K |
| 37 | Masters | Married-civ-spouse | Exec-managerial | 40 | United-States | <=50K |
| 49 | 9th | Married-spouse-absent | Other-service | 16 | Jamaica | <=50K |

Assume that the entries in columns ‘Age’ and ‘Hours per Week’ are of type `int` while the rest are of type `string`. In each part below, fill in the blanks of the Python expression to compute the described value. **You must use ONLY the line provided.** The code in the line should evaluate to the value described.

- (a) (2 pt) The proportion of adults who work exactly 40 hours a week

```
adults._____ / _____
adults.where('Hours per Week', 40).num_rows / adults.num_rows
```

- (b) (2 pt) The Relationship Status (string) of the youngest person in the table (you can assume there is only one youngest person)

```
adults._____
adults.sort('Age').column('Relationship Status').item(0)
```

- (c) (3 pt) An array consisting of the Education of the oldest adult (or adults, if more than one are oldest)

```
adults._____
adults.where('Age', max(adults.column('Age'))).column('Education')
```

- (d) (4 pt) A two-column table that has one row for each distinct country in the ‘Country’ column of the table `adults`, such that the first column ‘Country’ contains the name of the country, and the second column ‘Bachelors’ contains the average hours per week worked by adults in that country who have a Bachelors degree (and no higher education). For countries where no one has a Bachelors degree, the value of ‘Bachelors’ should be 0.

```
_____
adults.pivot('Education', 'Country', 'Hours per Week', 'np.mean').select(['Country', 'Bachelors'])
```

6. (7 points) Estimating a Percentile

Fill in the blanks below with code to define a function `ci_75` that constructs a confidence interval for the 75th percentile of a numerical population, as follows. The function takes the following arguments.

- **tbl**: A one-column table consisting of a random sample from the population; you can assume that the sample is large
- **reps**: a number of bootstrap repetitions; you can assume that users will enter a large integer

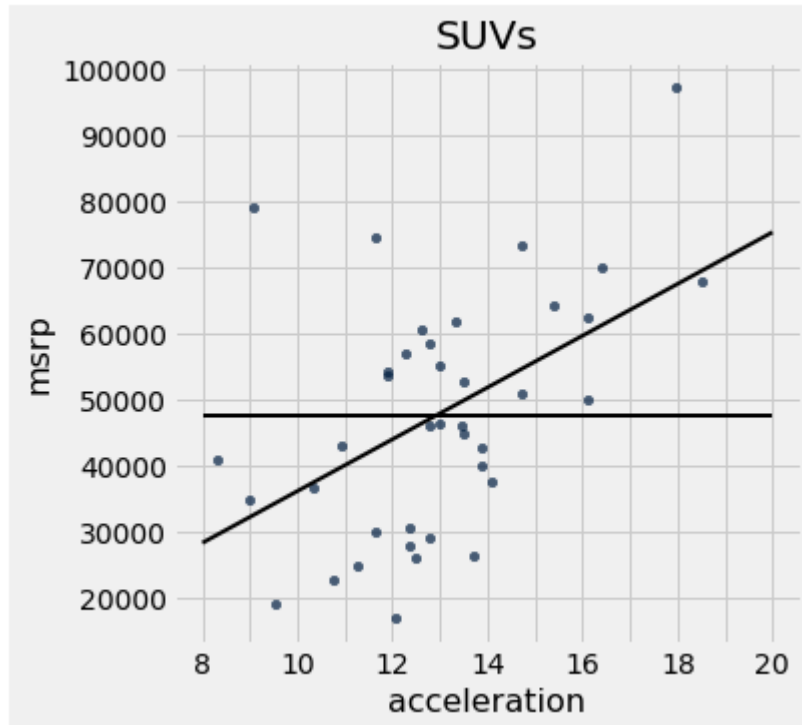
The function returns an array containing the endpoints of an approximate 95% bootstrap confidence interval for the 75th percentile of the population.

```
def ci_75(tbl, reps):  
  
    percentiles = make_array()  
  
    for i in np.arange(reps)  
  
        new_samp = tbl.sample()  
  
        new_percentile = percentile(75, new_samp.column(0))  
  
        percentiles = np.append(percentiles, new_percentile)  
  
    left_end = percentile(2.5, percentiles)  
  
    right_end = percentile(97.5, percentiles)  
  
    return make_array(left_end, right_end)
```

7. (13 points) Prediction and Error

In the scatter diagram below, each point represents a model of hybrid SUV. The variables are:

- acceleration, measured in kilometers per hour per second (it doesn't matter if you don't understand those units)
- msrp, an acronym for manufacturer's suggested retail price, in dollars



Some summary statistics:

- The average msrp is \$47,600 and the SD is \$18,000.
- The correlation between msrp and acceleration is 0.5.

The two straight lines:

- The flat line is at the level $y = \text{average msrp}$.
- The slanted line is the regression line for predicting msrp based on acceleration.

- (a) (2 pt) Fill in the blank and **explain**: If a hybrid SUV is one SD above average in acceleration, the regression prediction of its msrp is \$_____ above the average msrp.

Explanation:

The correlation between the two is .5, so as acceleration moves one SD above the average, we predict the msrp to be .5 SD above the average, which is equal to \$9000.

- (b) (3 pt) Pick one option and **explain**.

The average acceleration is closest to

- ☐ 12
- ☐ 12.5
- ☐ 13
- ☐ 13.5
- ☐ 14

Explanation:

13. The average of acceleration is the point of intersection between our two lines. Why? The regression line passes through the point (average of x, average of y), so whenever the regression line passes through the line (average of y), we know the corresponding x value is the average value of x. The largest mistake in this problem was noting that 13 is the intersection, but not saying why that meant it was the right answer still.

- (c) (2 pt) Pick one option to fill in the blank, and **explain**.

To predict msrp based on acceleration, Researcher Ave uses the flat line and Researcher Reg uses the regression line. The root mean squared error made by Researcher Ave is _____ the root mean squared error made by Researcher Reg.

- ☐ less than
- ☐ equal to
- ☐ greater than

Explanation:

greater than. Researcher Reg has the smallest possible rmse, as the regression line is the unique line which minimizes this. So it is definitely less than that of Researcher Ave.

- (d) (3 pt) Find the root mean squared error made by Researcher Reg, if it is possible to find it with the information given. **Explain your answer.** If it is not possible to find it, write NA and **explain your choice**.

Researcher Reg uses the regression line. The RMSE of the regression line's fitted points is just the standard deviation of the fitted points. This means that the rmse of researcher reg is $\sqrt{1 - .5^5} * 18000$.

- (e) (3 pt) Find the root mean squared error made by Researcher Ave, if it is possible to find it with the information given. **Explain your answer.** If it is not possible to find it, write NA and **explain your choice**.

Researcher Ave used the average line, which means we are looking for the RMSE around the average line. This is just the standard deviation! So our answer is 18000.

8. (6 points) Tests and Error

Researchers are conducting a test of hypotheses using 3% as the cutoff for the P-value.

(a) (3 pt) Pick the right option from among the following and **explain your choice**.

If the null hypothesis is true, the chance that the test reaches the correct conclusion is

- ☐ 100%
- ☐ 0%
- ☐ 2%
- ☐ 3%
- ☐ 97%
- ☐ 50%

Explanation:

Option 5. A 3% cut-off means, if the null is true, with probability of 3% we reject it incorrect. The other option is to remain consistent with the null, which means the test is correct. Since probabilities must add up to 1, this must have a 97% chance.

(b) (3 pt) Pick the right option from among the following and **explain your choice**.

The P-value of the test comes out to be 2%. The conclusion of the test is that

- ☐ the data support the null hypothesis more than they support the alternative, because there is a 98% chance that the null hypothesis is true.
- ☐ the data support the alternative hypothesis more than they support the null, because if the null were true then something unlikely has occurred.
- ☐ the data support the alternative more than they support the null, because there is only a 2% chance that the null hypothesis is true.

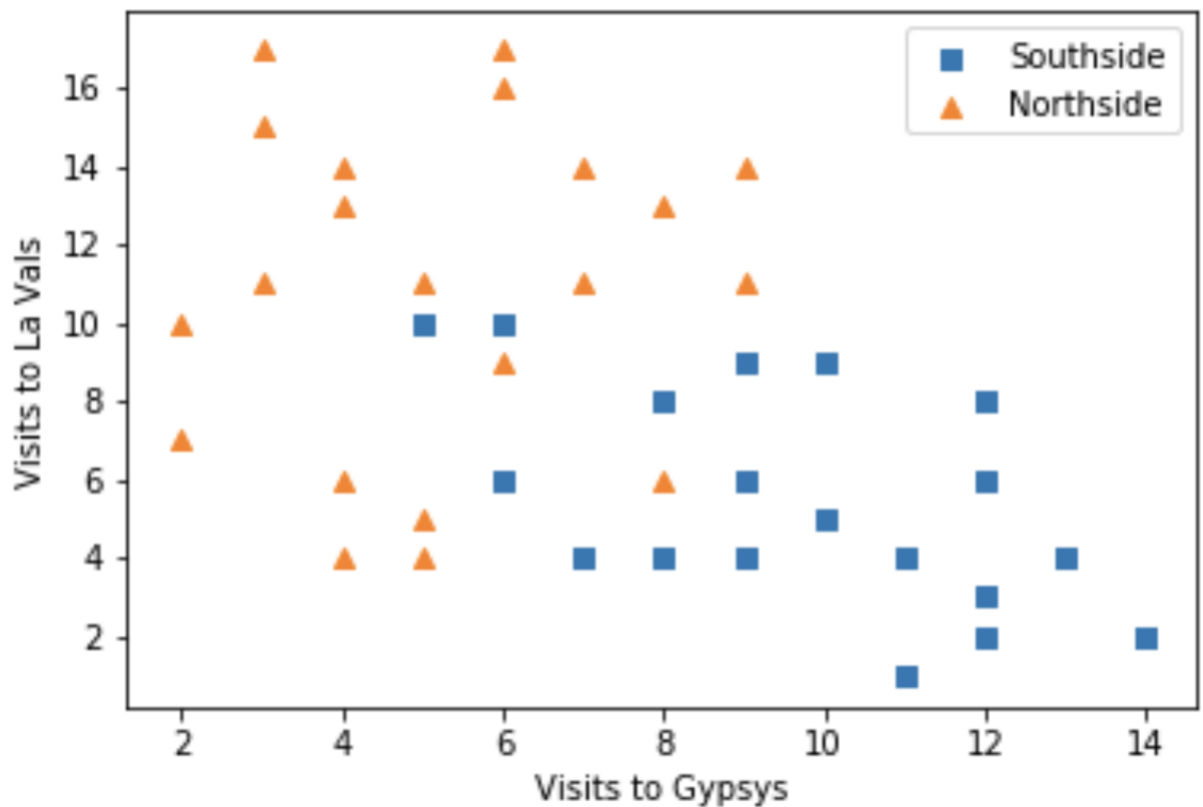
Explanation:

Option 2. It is not option 1 or 3; the null hypothesis is either true or not. There is no chance associated with it. Our p-value is less than our cut-off, so we say the data support the alternative hypothesis.

9. (7 points) Northside or Southside

A student is trying to build a classifier that classifies Berkeley students as residents of Northside or Southside. The student has a random sample of Berkeley students all of whom live on Northside or Southside. For each student she records whether the student lives on Northside or Southside, the number of times the student went to La Val's (on Northside) in the last 6 months, and the number of times the student went to Gypsy's (on Southside) in the last 6 months.

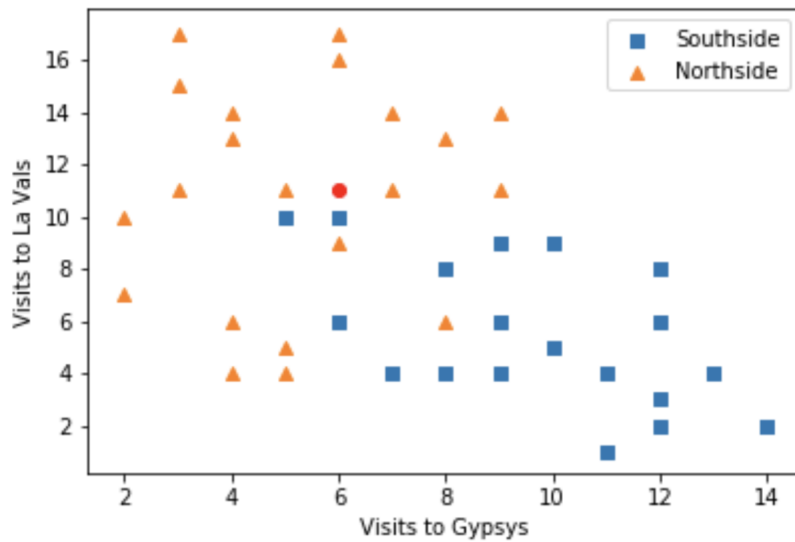
- (a) (3 pt) Draw an approximate 5-nearest-neighbors decision boundary on the scatter plot of the sample, shown below. No explanation is necessary.



Common Mistakes:

- We attempted to be very lenient in this problem, but if it was blatantly incorrect, we marked it off
- A decision boundary means a division where you can classify all points in a specific part of the division as one class
- Many students put a specific point on the graph and drew a decision boundary for it. This is not a correct decision boundary
- Decision boundaries are not lines which separate the classes in the graph either. In this decision boundary, it's possible for some Southside points to lie in a decision boundary where everything gets classified Northside (i.e. the two Southside points in the top left area)

The student is attempting to classify the point at (6, 11), represented by the circle in the graph below.



(b) (2 pt) Suppose she uses a 3-nearest-neighbors classifier. What will her classification be? **Explain.**

- ☐ Northside
☐ Southside

Explanation:

Northside: Three closest neighbors are one to the left, right, and below. Two of the three are Northside, so we classify the unknown point as Northside.

(c) (2 pt) Suppose she uses a 5-nearest-neighbors classifier instead. What will her classification be? **Explain.**

- ☐ Northside
☐ Southside

Explanation:

Northside: The five closest neighbors are evident on the graph. Three of those neighbors are Northside.

10. (7 points) Weighted K Nearest Neighbors

Instead of classifying a movie as ‘action’ or ‘romance’ according to the class of the majority of its k nearest neighbors in a training set, consider weighting each of the k nearest neighbors by its distance from the unclassified movie. The following method is proposed.

- Find the k nearest neighbors of the unclassified point.
- Let **Largest** be the largest of the distances of these k points from the unclassified point, and let **Smallest** be the smallest of the k distances.
- For each of the k nearest neighbors, define its weight as follows. Let d be its distance from the unclassified point. Calculate its weight as $(\text{Largest} - d) / (\text{Largest} - \text{Smallest})$
- Among the k nearest neighbors, find the total weight of the ‘action’ points as well as the total weight of the ‘romance’ points.
- Choose the class that has the larger total weight. If the two total weights are equal, choose whichever class you wish.

In order to classify an unclassified movie using the scheme above, complete the definition of a function `weighted_knn` that takes the following arguments:

- **tbl**: A two-column table in which the first column is labeled ‘Distance’ and the second column is labeled ‘Genre’. Each row of the table represents a movie in the training set. The first entry in the row is the distance of the movie from the unclassified movie. The second entry in the row is the genre (either ‘action’ or ‘romance’) of the movie.
- **k**: The value of k to use for k -nearest-neighbors

The function returns either ‘action’ or ‘romance’ using the weighted k -nearest-neighbors scheme proposed above.

```
def weighted_knn(tbl, k):

    nn = tbl.sort('Distance').take(np.arange(k))

    smallest = min(nn.column('Distance'))

    largest = max(nn.column('Distance'))

    weights = largest - nn.column('Distance') / (largest - smallest)

    weighted_table = nn.with_column('Weight', weights)

    by_genre = weighted_table.group('Genre', sum)
```

```
return by_genre.sort('Weights sum').column('Genre', descending=True).item(0)
```

11. (20 points) Voter Distributions

Before a Presidential election, Candidate A's campaign staff studied distributions of voter preferences by taking a random sample of voters in each of two states. Each sampled voter checked one of three boxes:

- Will vote for Candidate A
- Will not vote for Candidate A
- Undecided

The table `voters` contains one row for each sampled voter. The states are labeled State 1 and State 2. Here are the first four rows of `voters`, along with the output of `voters.group('State')`.

| State | Preference |
|---------|---------------------|
| State 1 | Will not vote for A |
| State 1 | Will vote for A |
| State 1 | Will vote for A |
| State 1 | Undecided |

| State | count |
|---------|-------|
| State 1 | 2000 |
| State 2 | 3000 |

- (a) (3 pt) The table `dists`, shown below, contains counts of sampled voters in different categories. For example, there are 540 Undecided voters in the sample from State 1.

| Preference | State 1 | State 2 |
|---------------------|---------|---------|
| Undecided | 540 | 600 |
| Will not vote for A | 660 | 1350 |
| Will vote for A | 800 | 1050 |

Fill in the blank so that the line of code produces the table `dists`.

```
dists = voters._____
```

```
dists = voters.pivot('State', 'Preference')
```

- (b) (4 pt) Are the distributions of preference different in the two states? State the null and alternative hypotheses that should be used to answer this question.

Null:

Alternative:

Null: The distributions of preference in the two states are the same. [The distributions in the samples are different due to chance.]

Alternative: The distributions of preference in the two states are different.

Common Mistakes:

- This question is focused on the overall distribution of voter preferences; not just whether or not someone will vote for candidate A
- A Null Hypothesis such as "The difference in distribution of preferences in the two states is due to chance" doesn't make sense. There is no chance associated with the distribution of preferences in the two states being the same. This would make sense if we replaced "in the two states with" "any sample".
- The alternative hypothesis should only claim that the distribution of preferences in the population of the two states is different. There's no claim about how they are different, just that they are.

- (c) (1 pt) What test statistic should be used to test the hypotheses in (b)?

Test Statistic example: TVD

Common Mistakes:

- We can not look at the absolute differences between the raw counts of the voter preferences. The sample sizes are different, so we must look at proportions

- (d) (1 pt) Complete the line of code below so that `obs_stat` is the observed value of the test statistic. You can use the table `dists` in your code.

```
obs_stat = _____
obs_stat = sum(abs(dists.column(1)/2000 - dists.column(2)/3000)) / 2
```

Common Mistakes:

- You must divide by 2000 and 3000 respectively, as this is how you get proportions (divide by the total number)
- If you used TVD as your test statistic from part c, you must remember to divide by 2

- (e) (3 pt) Explain how to simulate the test statistic under the null hypothesis. You will write the code in the next part.

Shuffle the column of voter preferences and match it to the original column of states. If the two distributions were the same, shuffling voting preferences would not make a difference as you treat their preferences equally.

Common Mistakes:

- In order to simulate the test statistic under the null, you must eventually calculate the test statistic. Failing to describe how to eventually calculate the test statistic and on what table to calculate it on was not considered a full explanation
- Remember that we shuffle, not bootstrap, the column of interest. That is why it is called a permutation test

- (f) (4 pt) Complete the definition of the function `simulate_stat` so that it returns one value of the test statistic simulated under the null hypothesis. The function takes no argument.

```
def simulate_stat()
    new_array = voters.sample(with_replacement=False).column('Preference')
    new_tbl = voters.select('State').with_column('Shuffled Prefs', shuffled)
    dists = new_tbl.pivot('State', 'Preference')
    return sum(abs(dists.column(1)/2000 - dists.column(2)/3000)) / 2
```

- (g) (4 pt) Fill in the code so that the last line evaluates to the P-value of the test based on 10,000 simulated values of the test statistic. Use `simulate_stat` and `obs_stat` in your code.

```
stats = make_array()
for i in np.arange(10000):
    stats = np.append(stats, simulate_stat())
np.count_nonzero(stats >= obs_stat) / 10000
```

Common Mistakes:

- `simulate_stat` is a function and hence must be called with parenthesis.
- Large values of the test statistic point to the alternative, so in order to calculate the p-value, we want to look at the simulated statistics that were greater than or equal to our observed test statistic