

NAME (FIRST, LAST): _____ SID: _____

TIME AND CONDITIONS: 45 minutes; closed book/notes/internet; no calculator/computer

QUESTIONS AND ANSWERS

- There are 8 questions, some of which have two parts.
- Not all questions will take the same amount of time. Some will go very fast.
- As in homework, some questions involve a short paragraph that you have to read. Reading quickly and accurately is highly recommended. It might help to mark key phrases as you read.
- You may answer any part of any question. If the answer to one part depends on another that you couldn't do, you can still provide an answer such as "The answer to part (a), divided by 2."
- When answers involve calculations that can't easily be done by mental arithmetic, we have asked you to leave the arithmetic unsimplified. Just leave them in any form that can be typed (perhaps laboriously) into a calculator to get the decimal answer.
- If a question asks for one line of code, you can take the space you need to write out the line by hand. You don't have to squeeze your handwritten code into one line of space on the exam page. But it should execute correctly if we type it into one line in a notebook cell.
- Explanations are expected to be concise. One or two clear sentences should be enough. Calculations and code are sufficient as explanations.

GRADING

- The exam is worth 50 points.
- You get 1 point for turning in a test with your name and SID on it, and 1 more point if you write each of your answers within the space provided. See **FORMAT** below.
- Each of the 8 questions is worth 6 points. If a question has two parts, each part is worth 3 points.
- We will give partial credit, but only for substantial progress towards a correct answer. We get to decide what "substantial progress" means.
- Commit yourself to a single answer for each part of each question. If you give multiple answers (such as both True and False), please don't expect credit, even if the right answer is among those that you gave.

FORMAT

- There is space for your answer below each question. **Please do not write outside the black boundary;** the scanner and Gradescope won't read it.
- If you need scratch paper, please use the backs of the pages of the exam, but be aware that they will not be graded.
- A reference sheet of class content will be provided. But it does not contain everything that was covered in class.

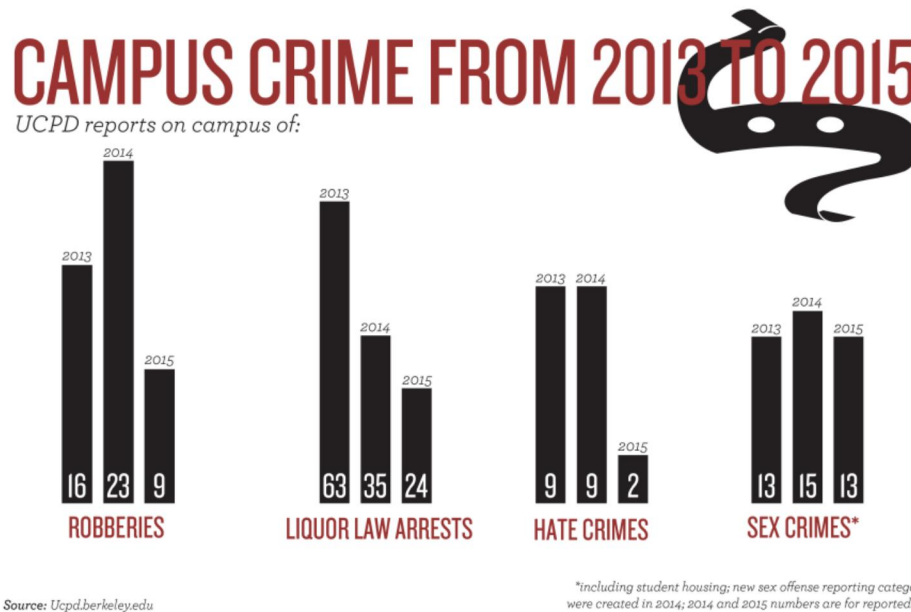
HONOR CODE

Data Science and the entire academic enterprise are based on one quality – integrity. We are all part of a community that doesn't fabricate evidence, doesn't fudge data, doesn't steal other people's work, doesn't lie and cheat. You trust that we will treat you fairly and with respect. We trust that you will treat us and your fellow students fairly and with respect. **Please abide by UC Berkeley's Honor Code:**

"As a member of the UC Berkeley community, I act with honesty, integrity, and respect for others."

Name: _____

1. The graphic below appeared on the front page of the Daily Cal on Tuesday of last week.



(a) One important choice made by the graphic designers makes it difficult to visually compare the numbers of crimes in various categories. What was that choice, and why does it make visual comparisons difficult?

(b) The data have been entered into a table called **crime**, shown below.

Year	Robberies	Liquor Law Arrests	Hate Crimes	Sex Crimes
2013	16	63	9	13
2014	23	35	9	15
2015	9	24	2	13

Use the table **crime** in one line of Python code that produces a graphic that corrects the problem in (a), and **explain** why the problem will be fixed.

Name: _____

2. The table **colors** is created as follows:

```
colors = Table().with_column(  
    'Color', make_array('blue', 'green', 'red', 'purple')  
)
```

Find a numerical expression for the chance that all the entries of **colors.sample(3).column(0)** are the same. Don't simplify the arithmetic.

Name: _____

3. Fill in each blank with the best choice from the table of items below. **No explanations are necessary.**

You can use items more than once. Don't write out the words or phrases you picked; just enter the item number in the answer table.

1.	approximate	2.	empirical
3.	sample	4.	population
5.	random sample	6.	estimate
7.	parameter	8.	statistic
9.	low	10.	high
11.	bias	12.	confounding
13.	likely	14.	unlikely
15.	variability	16.	distribution
17.	approximation	18.	probability
19.	convenience	20.	control

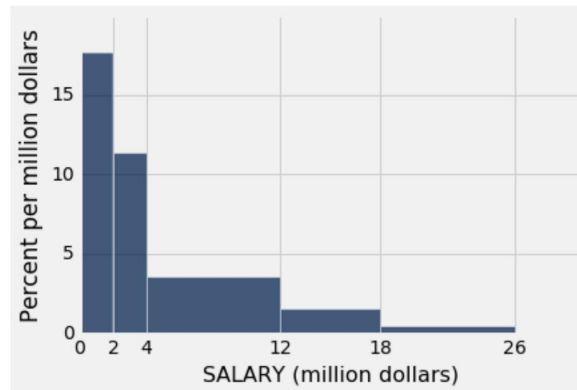
- Important criteria for a statistic to be a good estimate of a population ____ I ____ include ____ II ____ bias and low ____ III ____.
- The ____ IV ____ distribution of a large ____ V ____ is likely to resemble the ____ VI ____ of the population.

Answer Table

Blank	Item Number
I	
II	
III	
IV	
V	
VI	

Name: _____

4. The table **nba** has a column labeled **SALARY** containing the 2015-2016 salaries of NBA players. Here is the output of `nba.select('SALARY').hist(bins = make_array(0, 2, 4, 12, 18, 26))` along with the heights of the bars.



bin (million dollars)	[0, 2)	[2, 4)	[4, 12)	[12, 18)	[18, 26)
height (percent per million dollars)	17.63	11.39	3.60	1.60	0.45

The interval $[a, b)$ contains all values that are greater than or equal to a and less than b .

(a) Which bin contains more players: $[2, 4)$ or $[4, 12)$? **Explain** your choice.

b) To see some more detail in the $[4, 12)$ range, the histogram will be redrawn with bins as shown below. The display includes the heights that are available from above.

bin (million dollars)	[0, 2)	[2, 4)	[4, 6)	[6, 12)	[12, 18)	[18, 26)
height (percent per million dollars)	17.63	11.39	(i)	(ii)	1.60	0.45

The expression `nba.num_rows` evaluates to 417.

The expression `nba.where('SALARY', are.between(4, 6)).num_rows` evaluates to 56.

If possible, provide a numerical expression for each missing height (do not simplify the arithmetic). If this is not possible, explain why not.

(i)

(ii)

Name: _____

5. In a Public Health Service study of the effect of smoking, participants were separated into groups according to gender and age. In each subgroup, the observation was the same – those who had never smoked were on average somewhat more healthy than the current smokers, but the current smokers were on average quite a bit more healthy than those who had recently stopped smoking.

Based on these results, should the current smokers have been advised not to give up smoking? If not, provide at least one factor that helps explain the observation.

Name: _____

6. A clinic is trying to run a randomized controlled trial of a new appointment system. The clinic is open Monday through Saturday each week. As a way of randomization, all appointment requests that come in on Mondays, Wednesdays, and Fridays are assigned to the new system. All appointment requests that come in on Tuesdays, Thursdays, and Saturdays are assigned to the old system.

At the end of a few weeks, 212 appointment requests have been assigned to the new system and 277 to the old one.

Clinician *A* says, “This method is biased against the people who call on on Mondays, Wednesdays, and Fridays. If we had tossed a fair coin for each request, we wouldn’t have ended up with so few requests assigned to the new system.”

Clinician *B* says, “No, the results are just like tossing a fair coin.”

(a) Provide null and alternative hypotheses that reflect the views of the two clinicians.

Null hypothesis:

Alternative hypothesis:

(b) A test of the hypotheses in part (a) will be performed. The test statistic is the **number of requests assigned to the new system**.

Fill in the blanks in the code below so that the output of the final line is the empirical P-value of the test. Make sure your code is consistent with the hypotheses in part (a).

Two arrays have been created for you. You are welcome to use their names in your code if you wish.

- The array `array2` consists of two entries, one of which is the string ‘New’ and the other is ‘Old’.
- The array `array489` consists of 489 entries, 212 of which are ‘New’ and 277 are ‘Old’.

Here is the code for you to complete.

```
stats = make_array()
```

```
for i in np.arange(5000):
```

```
    draws = np.random.choice(_____, _____)
```

```
    number_new = np.count_nonzero(_____)
```

```
    stats = np.append(_____, number_new)
```

```
empirical_P = np._____
```

```
empirical_P
```

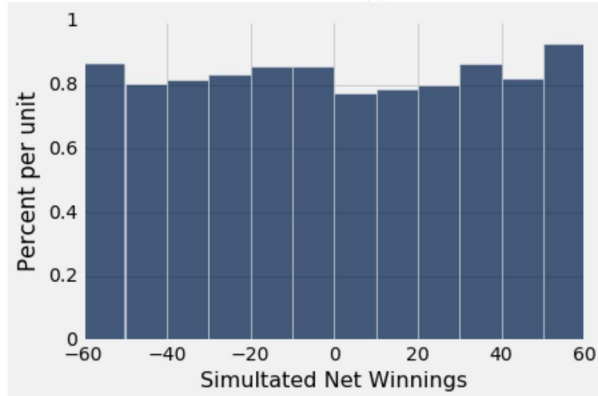
Name: _____

7. My friend and I play the following gambling game repeatedly:

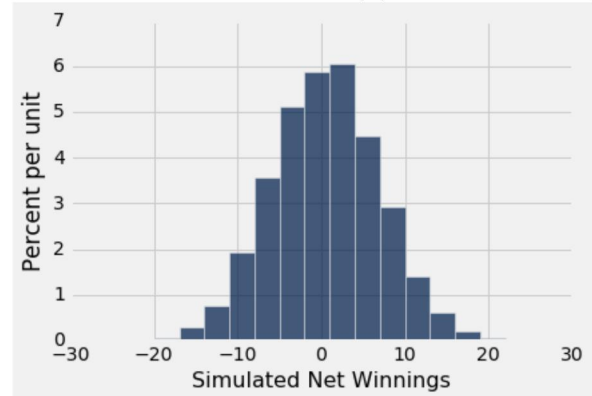
A fair die is rolled. If it shows one or six spots, I pay my friend \$1. If it shows two or five spots, my friend pays me \$1. If it shows three or four spots, neither of us pays any money.

One of the histograms below is an empirical histogram of my overall net winnings in 60 games, based on 10,000 repetitions of finding the overall net gain in 60 games. **Which histogram is it, and why?**

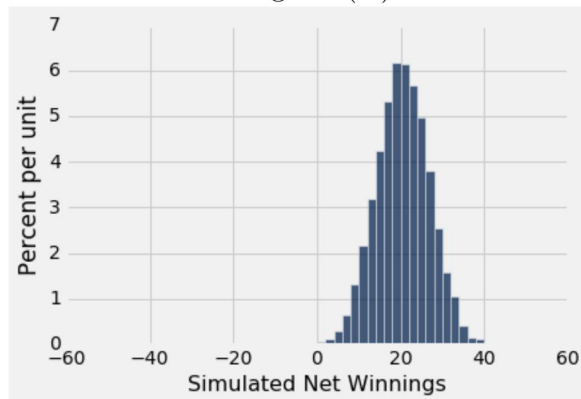
Histogram (i)



Histogram (ii)



Histogram (iii)



Histogram (iv)



Justification:

Name: _____

8. The table **voters** contains one row for each voter in a large national sample. There are two columns:

Column 0: **State** is the state where the voter is registered. The states are denoted by their two-letter abbreviations, so that California is the string '**CA**' and Florida is '**FL**'.

Column 1: **Choice** has three categories: '**Candidate A**', '**Candidate B**', and '**Undecided**'.

Write code to calculate the total variation distance between the **Choice** distributions of California and Florida. You can assume that the table contains at least one voter in each **Choice** category for both the states.

Use the following steps.

A table of the data for all the California voters, containing just one column: **Choice**

ca_voters = **voters**._____

An array containing the **Choice** distribution (in proportions) for California

ca_dist = **ca_voters**._____

Now assume you have created **fl_voters** and **fl_dist** correspondingly for Florida. There is no need to write out the code for those.

Find the total variation distance:

tvd = _____