

## Problem Set #5

### Due: 1:00pm on Monday, Aug 3

With problems by Mehran Sahami, Chris Piech, David Varodayan, Anand Shankar, and Lisa Yan

**For each problem, briefly explain/justify how you obtained your answer.** In fact, most of the credit for each problem will be given for the derivation/model used as opposed to the final answer.

### Written Problems

**For each problem, briefly explain/justify how you obtained your answer.** In fact, most of the credit for each problem will be given for the derivation/model used as opposed to the final answer.

- The joint probability density function of continuous random variables  $X$  and  $Y$  is given by:

$$f_{X,Y}(x, y) = c \frac{y}{x} \quad \text{where } 0 < y < x < 1$$

- What is the value of  $c$  in order for  $f_{X,Y}(x, y)$  to be a valid probability density function?
  - Are  $X$  and  $Y$  independent? Explain why or why not.
  - What is the marginal density function of  $X$ ?
  - What is the marginal density function of  $Y$ ?
  - What is  $E[X]$ ?
- A robot is located at the *center* of a square world that is 10 kilometers on each side. A package is dropped off in the robot's world at a point  $(x, y)$  that is uniformly (continuously) distributed in the square. If the robot's starting location is designated to be  $(0, 0)$  and the robot can only move up/down/left/right parallel to the sides of the square, the distance the robot must travel to get to the package at point  $(x, y)$  is  $|x| + |y|$ . Let  $D$  = the distance the robot travels to get to the package. Compute  $E[D]$ .
  - Let  $X$ ,  $Y$ , and  $Z$  be independent random variables, where  $X \sim \mathcal{N}(\mu_1, \sigma_1^2)$ ,  $Y \sim \mathcal{N}(\mu_2, \sigma_2^2)$ , and  $Z \sim \mathcal{N}(\mu_3, \sigma_3^2)$ .
    - Let  $A = X + Y$ . What is the distribution (along with parameter values) of  $A$ ?
    - Let  $B = 4X + 3$ . What is the *joint* distribution (along with parameter values) of  $B$  and  $Z$ ? (Hint: Bivariate Normal)
    - Let  $C = aX - b^2Y + cZ$ , where  $a$ ,  $b$ , and  $c$  are real-valued constants. What is the distribution (along with parameter values) of  $C$ ? Show how you derived your answer.

## Coding/Written Problems

Here are some rules that apply to **all the coding questions**:

- For coding parts, implement the corresponding functions of the provided Python3 files, available for download from the course website. For written parts, write your answer in the PDF that gets submitted to Gradescope. Do not rename any files.
  - Your code will be autograded.
  - Do not use global variables.
  - You may define helper functions if you wish.
  - All data files are located in the `datasets` directory, so you can use the provided helper function `get_filepath(filename)` to resolve paths, e.g., `get_filepath(email.txt)` returns `datasets/email.txt`. This function is included with each file, other than the extra credit article submission file.
6. Did you know that computers can identify you not only by what you write, but also by how you write? Coursera uses Biometric Keystroke signatures for plagiarism detection. If you cannot write a sentence with the same statistical distribution of key press timings as in your previous work, they assume that you are not the person sitting behind the computer. In this problem we provide you with several data files:
- `personKeyTimingA.txt` has keystroke timing information for a user A writing a passage. The first column is the time in milliseconds (since the start of writing) when the user hit each key. The second column is the key that the user hit.
  - `personKeyTimingB.txt` has keystroke timing information for a second user (user B) writing the same passage as the user A. Even though the content of the passage is the same the timing of how the second user wrote the passage is different.
  - `email.txt` has keystroke timing information for an unknown user. We would like to know if the author of the email was user A or user B.

Let  $X$  and  $Y$  be random variables for the duration of time, in milliseconds, for users A and B (respectively) to type a key. Assume that each keystroke from a user has a duration that is an independent random variable with the same distribution.

For coding problems, write your answer in the relevant function of `cs109_pset5_email.py`. Follow the same coding guidelines as listed above. For written problems, write your answer in the PDF that gets submitted to Gradescope.

- a. **[Coding]** Complete the function `part_a` provided in the starter code. This function takes in a `filename` (which is either `personKeyTimingA.txt` or `personKeyTimingB.txt`) and should return  $E[X]$  or  $E[Y]$ , respectively.
- b. **[Coding]** Complete the function `part_b` provided in the starter code. This function should return  $E[X^2]$  or  $E[Y^2]$ , depending on which file is supplied as `filename`.

- c. **[Written]** Use your answers to part (a) and (b) and approximate  $X$  and  $Y$  as Normal random variables with mean and variance that match their biometric data. Report both distributions.
- d. **[Written]** Calculate the ratio of the probability that user A wrote the email over the probability that user B wrote the email. You do not need to submit code, but you should include the formula that you attempted to calculate and a short description (a few sentences) of how your code works.
7. **[Coding + Written]** Stanford’s HCI class runs a massive online class that was taken by ten thousand students. The class used peer assessment to evaluate students’ work. We are going to use their data to learn more about peer graders. In the class, each student has their work evaluated by 5 peers and every student is asked to evaluate 6 assignments: five peers and the “control assignment” (the graders were unaware of which assignment was the control). All 10,000 students evaluated the same control assignment, and the scores they gave are in the file `peerGrades.csv`. You may use simulations to solve any part of this question.

For coding problems, write your answer in the relevant function of `cs109_pset5_hci.py`. Follow the same coding guidelines as listed above. For written problems, write your answer in the PDF that gets submitted to Gradescope.

- a. **[Coding]** What is the sample mean of the 10,000 grades to the control assignment? Implement the `part_a` function, which should return this quantity as a float.

For parts (b) and (c), you’ll need to run some simulations. To get credit from the autograder, you’re **required** to abide by the following guidelines:

- Run the algorithm for exactly 10,000 iterations.
  - You’ll need to draw random samples with replacement from an array of grades. To do so, you must use the `np.random.choice` function, which you can call like so:  
`sample = np.random.choice(name-of-array, size-of-random-sample, replace=True)`. Do not use any other function to generate random samples.
  - Use the `np.mean`, `np.median`, and `np.var` functions to calculate the mean, median, and variance of a list or numpy array.
- b. **[Coding]** Students could be given a final score which is the *mean* of the 5 grades given by their peers. Imagine the control experiment had only received 5 peer-grades. What is the variance of the mean grade that the control experiment would have been given? Implement the `part_b` function, which should return this quantity as a float.
- c. **[Coding]** Students could be given a final score which is the *median* of the 5 grades given by their peers. Suppose the control experiment had only received 5 peer-grades. What is the variance of the median grade that the control experiment would have been given? Implement the `part_c` function, which should return this quantity as a float.
- d. **[Written]** Would you use the mean or the median of 5 peer grades to assign scores in the online version of Stanford’s HCI class? Hint: it might help to visualize the scores. Feel free to write code to help you answer this question, but for this question we’ll solely evaluate your written answer in the PDF that you upload to Gradescope.

8. **[Coding + Written]** In this problem you are going to learn how to use and misuse  $p$ -values for experiments that are called *A/B tests*. These experiments are ubiquitous. They are a staple of both scientific experiments and user interaction design.

Suppose you are working at Coursera on new ways of teaching a concept in probability. You have two different learning activities `activity1` and `activity2` and you want to figure out which activity leads to better learning outcomes. Over a two-week period, you randomly assign each student to be given either `activity1` or `activity2`. You then evaluate each student's learning outcomes by asking them to solve a set of problems.

The data (the activity shown to each student and their measured learning outcomes) are found in the file `learningOutcomes.csv`. For coding problems, write your answer in the relevant function of `cs109_pset5_coursera.py`. Follow the same coding guidelines as listed above. For written problems, write your answer in the PDF that gets submitted to Gradescope.

- a. **[Coding]** What is the difference in sample means of learning outcomes between students who were given `activity1` and students who were given `activity2`? Write your answer in the `part_a` function, which should return a float (i.e. the difference in sample means).
- b. **[Coding]** Write code to estimate the  $p$ -value (using the bootstrap method) for the observed difference in means reported in part (a). In other words: assuming that the learning outcomes for students who had been given `activity1` and `activity2` were identically distributed, what is the probability that you could have sampled two groups of students such that you could have observed a difference of means as extreme, or more extreme, than the one calculated from your data in part (a)? Write your answer in the `part_b` function, which should return a float. Here are some guidelines to follow:
  - Just like in the previous problem, you are **required** to use the `np.random.choice` method with `replace=True` to generate random samples.
  - For the bootstrap algorithm, you should use 10,000 iterations, i.e. you should resample 10,000 times.
  - If you have two lists `a` and `b`, you can create a new list containing all the elements of `a` followed by all the elements of `b` by writing `a + b`

Scientific journals have traditionally accepted an experiment's result as "statistically significant" if the  $p$ -value is below 0.05. By definition, this standard means that 5% of findings published in these journals are in fact not true, but just false positives. The scientific community is beginning to move away from using arbitrary  $p$ -value thresholds to determine whether a result is publishable. For example, see this 2019 editorial in the journal *Nature*: <https://www.nature.com/articles/d41586-019-00874-8>.

You are now troubled by the  $p$ -value you obtained in part (b), so you decide to delve deeper. You investigate whether learning outcomes differed based on the background experience of students. The file `background.csv` stores the background of each student as one of three labels: more experience, average experience, less experience.

- c. **[Written]** For each of the three backgrounds, calculate a difference in means in learning outcome between `activity1` and `activity2`, and the  $p$ -value of that difference.

You'll almost certainly need to write code in this question, and we've provided an `optional_function` that you can use, which gets called by our provided main method. However, we won't grade any code for this part. We'll only grade what you include in your answer PDF.

- d. **[Written]** Your manager at Coursera is concerned that you have been “*p*-hacking,” which is also known as data dredging: [https://en.wikipedia.org/wiki/Data\\_dredging](https://en.wikipedia.org/wiki/Data_dredging). In one sentence, explain why your results in part (c) are not the result of *p*-hacking.