

Lisa Yan and Jerry Cain
CS 109

Quiz #2
November 2, 2020

CS109 Quiz #2 Solutions

Take-Home Quiz information

Exam # 00000

Each quiz will be a 47-hour open-book, open-note exam. We have designed this quiz to approximate about 1-3 hours of active work (*before* typesetting).

- You can submit multiple times; we will only grade the last submission you submit before 1:00pm (Pacific time) on Friday, October 30th. No exam submissions will be accepted late. When uploading, please assign pages to each question. Failure to do so will result in a 2-point deduction. ***Please double-check that you submit the right file.***
- You should upload your submission as a PDF to Gradescope. We provide a LaTeX template if you find it useful, but we will accept any legible submission. You may also find the CS109 Probability LaTeX reference useful: <https://www.overleaf.com/read/wyhtzmdsfwkb>
- Course staff assistance will be limited to clarifying questions of the kind that might be allowed on a traditional, in-person exam. If you have questions during the exam, please ask them as private posts via our discussion forum. We will not have any office hours for answering quiz questions during the quiz.
- **For each problem, briefly explain/justify how you obtained your answer** at a level such that a future CS109 student would be able to understand how to solve the problem. It is fine for your answers to be a well-defined mathematical expression including summations (but not integrals), products, factorials, exponentials, and combinations, unless the question *specifically* asks for a numeric quantity or closed form. Where numeric answers are required, fractions are fine.

Honor Code Guidelines for Take-Home Quizzes

This exam must be completed individually. It is a violation of the Stanford Honor Code to communicate with any other humans about this exam (other than CS109 course staff), to solicit solutions to this exam, or to share your solutions with others.

The take-home exams are open-book: open lecture notes, handouts, textbooks, course lecture videos, and internet searches for conceptual information (e.g., Wikipedia). Consultation of other humans in any form or medium (e.g., communicating with classmates, asking questions on sites like Chegg or Stack Overflow) is prohibited. All work done with the assistance of any external material in any way (other than provided CS109 course materials) must include citation (e.g., “Referred to Wikipedia page on X for Question 2.”). Copying solutions is unacceptable, even with citation. If by chance you encounter solutions to the problem, navigate away from that page before you feel tempted to copy.

If you become aware of any Honor Code violations by any student in the class, your commitments under the Stanford Honor Code obligate you to inform course staff. ***Please remember that there is no reason to violate your conscience to complete a take-home exam in CS109.***

I acknowledge and accept the letter and spirit of the Honor Code:

Name (typed or written): _____

1 Undergraduates and Part-Time Jobs [22 points]

According to the California Center for Overextended Students, 74% of Stanford undergraduates and 86% of Berkeley undergraduates maintain part-time jobs in addition to their coursework, all independently of each other.

- a. (4 points) You interview five Stanford undergraduates and five Berkeley undergraduates. What is the probability that exactly nine of the ten hold part-time jobs?

Answer. Let S equal the number of Stanford students who have a full-time job and B equal the number of Berkeley undergraduates who have full-time jobs. We want to compute $P(S + B = 9)$, which is trivially a convolution, but more easily expressed as a sum of two mutually exclusive events.

$$\begin{aligned} P(S + B = 9) &= P(S = 5, B = 4) + P(S = 4, B = 5) \\ &= P(S = 5) P(B = 4) + P(S = 4) \cdot P(B = 5) \\ &= \binom{5}{4} (0.74)^4 (0.26) \binom{5}{5} (0.86)^5 + \binom{5}{5} (0.74)^5 \binom{5}{4} (0.86)^4 (0.09) \\ &= \binom{5}{4} \binom{5}{5} (0.74)^4 (0.26) (0.91)^5 + \binom{5}{5} \binom{5}{4} (0.74)^5 (0.86)^4 (0.09) \\ &= 5 \cdot (0.74)^4 (0.26) (0.86)^5 + 5 \cdot (0.74)^5 (0.86)^4 (0.14) \\ &= 0.27692288099999995 \end{aligned}$$

- b. (5 points) Assume you interview Stanford undergraduates, one after another, to survey who holds a part-time job. Let S_k be the number of Stanford undergraduate that must be interviewed until you find the k^{th} student who holds a part-time job. What is $P(S_{14} = 14)$?

Answer. The probability that we need to interview 14 Stanford undergraduates before we finally find 14 students who hold part-time jobs is given by:

$$\begin{aligned} P(S_{14} = 14) &= \binom{14-1}{14-1} \cdot (0.74)^{14} \cdot (0.26)^{14-14} \\ &= 0.27692288099999995 \end{aligned}$$

- c. (7 points) Suppose there are 7000 Stanford undergraduates and 33000 Berkeley undergraduates, and you choose two students at random from the 40000 combined and learn they each have part-time jobs? What is the conditional probability that both students attend Berkeley, given both have part-time jobs?

Answer. Let BB represent the event that both undergraduates attend Berkeley, and let WW represent the event that both undergraduates work. First, we compute the probability that both undergraduates attend Berkeley, which is more straightforward than the probabilities we

need to compute: $P(BB) = (\frac{33000}{40000})^2$. The probability that both undergraduates work given they each attend Berkeley is equally as straightforward: $P(WW|BB) = 0.86^2$. Now, computing the probability that both randomly chosen students work is only slightly more involved: $P(WW) = P(WW|BB) P(BB) + P(WW|BS) P(BS) + P(WW|SB) P(SB) + P(WW|SS) P(SS)$. This is involved enough that we should derive it:

$$\begin{aligned} P(WW) &= P(WW|BB) P(BB) + P(WW|BS) P(BS) + P(WW|SB) P(SB) + P(WW|SS) P(SS) \\ &= P(WW|BB) P(BB) + 2 P(WW|BS) P(BS) + P(WW|SS) P(SS) \\ &= (\frac{33000}{40000})^2 \cdot 0.86^2 + 2 \cdot (\frac{7000}{40000})(\frac{33000}{40000}) \cdot 0.86 \cdot 0.74 + (\frac{7000}{40000})^2 \cdot 0.74^2 \end{aligned}$$

This, Stanford students, is classic Bayes Theorem.

$$\begin{aligned} P(BB|WW) \cdot P(WW) &= P(WW|BB) \cdot P(BB) \\ P(BB|WW) &= \frac{P(WW|BB) \cdot P(BB)}{P(WW)} \\ &= 0.27692288099999995 \end{aligned}$$

- d. (6 points) Suppose you interview 120 Stanford undergraduates individually. What is the approximate probability that strictly more than 90 of them hold part-time jobs?

Answer. Because 120 is large and the variance (which is 300) is greater than 20, we can approximate $\text{Bin}(120, 0.74)$ the distribution as Normal, with μ equal to 450 and σ^2 of 300. We need to reframe the our 90 in terms of the number of standard deviations below the estimated mean.

Let X represent the number of Stanford undergraduates who hold part-time jobs, and let Y approximate X using $N(450, 300)$. That means that:

$$\begin{aligned} P(X > 90) &= P(Y > 90.5) \\ &= 1 - P(Y < 90.5) \\ &= 1 - \Phi\left(\frac{90.5 - 450}{17}\right) \\ &= 1 - \Phi(-1.45) \\ &= \Phi(1.45) \\ &= 0.27692288099999995 \end{aligned}$$

2 Fish Sticks [12 points]

Fish Sticks (a frozen meal company) wants to analyze the success of their new home-delivery website. Suppose users independently visit the homepage at an average rate of 3 users per minute.

- a. (2 points) What is the expected number of users who visit the homepage in the next 5 minutes?

Answer. Let X be the number of visitors in 5 minutes. $X \sim \text{Poi}(\lambda = 15)$, where using stoichiometry, the average number of (independent) visitors in 5 minutes is 15 users per minute. $E[X] = \lambda = 15$.

Alternatively, let X_i be the number of visitors in the i^{th} minute, where $X_i \sim \text{Poi}(\lambda = 3)$ and $E[X_i] = 3$. Note that $X = X_1 + X_2 + \cdots + X_5$ and therefore by linearity of expectation, $E[X] = \sum_{i=1}^5 E[X_i] = 5(3) = 15$.

- b. (4 points) What is the probability that there are exactly 20 users who visit the homepage in the next 5 minutes?

Answer. Using our definition of X from above, $X \sim \text{Poi}(\lambda = 15)$. Therefore

$$P(X = 20) = \frac{15^{20}e^{-15}}{20!} \approx 0.0418.$$

Each user that visits the homepage downloads exactly 5 MB of website data. We define **homepage load** in a minute to be the total amount of website data downloaded across all users who visit the homepage in that minute, e.g., 10 homepage visitors in the next minute means the homepage load in that minute is 50 MB. As before, users still independently visit the homepage at an average rate of 3 users per minute.

- c. (2 points) What is the expected homepage load in the next minute?

Answer. Let D be the total data downloaded (in MB) across all users who visit the homepage in a minute. If V is the number of visitors in that minute (where $V \sim \text{Poi}(\lambda = 3)$), then $D = 5V$. Therefore by linearity of expectation, $E[D] = E[5V] = 5E[V] = 5(3) = 15$.

- d. (4 points) What is the probability that the homepage load is exactly 20 MB in the next minute?

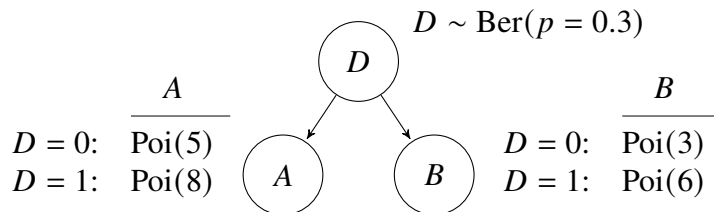
Answer. Using the definitions of D and V as in part (c), we are looking to compute $P(D = 20) = P(5V = 20) = P(V = 4)$. Since $V \sim \text{Poi}(3)$,

$$P(D = 20) = P(V = 4) = \frac{3^4e^{-3}}{4!} \approx 0.1680.$$

Note that we cannot model D as a Poisson random variable, because D can only take on non-negative values that are multiples of 5 and not the entire range of support of a typical Poisson (non-negative integers).

3 Stanford Fish Sticks [30 points]

Fish Sticks (the same frozen meal company) now wants to model their hourly homepage traffic from Stanford. The company decides to model two different behaviors for homepage visits according to the Bayesian Network on the right:



A and B are the numbers of Stanford students and faculty, respectively, who visit the Fish Sticks homepage in an hour. Since Fish Sticks does not know when Stanford people eat, the company models demand as a “hidden” Bernoulli random variable D , which determines the distribution of A and B . Recall that in a Bayesian Network, random variables are conditionally independent given their parents. For example, given $D = 0$, $A \sim \text{Poi}(5)$ and $B \sim \text{Poi}(3)$, two independent random variables.

- a. (6 points) Given that 6 users from group A visit the homepage in the next hour, what is the probability that $D = 0$?

Answer. Note that given $D = 0$, $A \sim \text{Poi}(\lambda = 5)$, and given $D = 1$, $A \sim \text{Poi}(\lambda = 8)$. By Bayes’ Theorem,

$$\begin{aligned}
 P(D = 0|A = 6) &= \frac{P(A = 6|D = 0)P(D = 0)}{P(A = 6|D = 0)P(D = 0) + P(A = 6|D = 1)P(D = 1)} \\
 &= \frac{\frac{5^6 e^{-5}}{6!} (1 - 0.3)}{\frac{5^6 e^{-5}}{6!} (1 - 0.3) + \frac{8^6 e^{-8}}{6!} (0.3)} \\
 &= \frac{5^6 e^{-5} (1 - 0.3)}{5^6 e^{-5} (1 - 0.3) + 8^6 e^{-8} (0.3)} \approx 0.7364
 \end{aligned}$$

- b. (10 points) What is the probability that in the next hour, the *total* number of users who visit the homepage from groups A and B is equal to 12, i.e., what is $P(A + B = 12)$?

Answer. By Law of Total Probability,

$$P(A + B = 12) = P(A + B = 12|D = 0)P(D = 0) + P(A + B = 12|D = 1)P(D = 1).$$

A and B are conditionally independent Poisson random variables given D , and therefore $A + B|D = 0 \sim \text{Poi}(\lambda = 8)$ and $A + B|D = 1 \sim \text{Poi}(\lambda = 14)$. Using the Poisson PMF,

$$P(A + B = 12) = \frac{8^{12} e^{-8}}{12!} \cdot (1 - 0.3) + \frac{14^{12} e^{-14}}{12!} \cdot (0.3) \approx 0.0632.$$

- c. (14 points) Simulate $P(A + B = \text{total})$, where `total = 12`, by implementing the `infer_prob_total(total, ntrials)` function below using rejection sampling.
- `total` is the total number of users from groups A and B in the event $A + B = \text{total}$.
 - `ntrials` is the number of observations to generate for rejection sampling.

- `prob` is the return value to the function, where $\text{prob} \approx P(A + B = \text{total})$.
- The function call is implemented for you at the bottom of the code block.

You can call the following functions from the `scipy` package:

- `stats.bernoulli.rvs(p)`, which randomly generates a 0 or 1 with probability p
- `stats.poisson.rvs(λ)`, which randomly generates a value according to a Poisson distribution with parameter λ

You are not required to use lists or NumPy arrays in this question (but you can if you want).

Pseudocode is fine as long as your code accurately conveys your approach. We are not grading on style nor syntax.

```
import numpy as np
from scipy import stats

def infer_prob_total(total, ntrials):
    n_samples_event = 0
    for i in range(ntrials):
        d = stats.bernoulli.rvs(0.3)
        user_sum = 0
        if d == 0:
            user_sum += stats.poisson.rvs(5) + stats.poisson.rvs(3)
        else:
            user_sum += stats.poisson.rvs(8) + stats.poisson.rvs(6)

        if user_sum == 12:
            n_samples_event += 1

    prob = n_samples_event/ntrials
    return prob

ntrials = 50000
total = 12
print("Simulated P(A+B)=", infer_prob_total(total, ntrials))
```

4 Mutually Recursive Code Analysis [16 points]

After much research, you’ve finally arrived at function to how many pairs of flip flops you should bring to the beach, since flip flops tend to disappear in the sand and get washed away in the tide, and you need to have at least one pair when you leave.

Consider the following implementation of the mutually recursive flip and flop functions:

```
from scipy import stats
from random import choice

def flip():
    eel = choice([1, 2, 3])
    if eel == 1: return 3
    if eel == 2: return 5 + flip()
    return 11 + flop()

def flop():
    fish = stats.poisson.rvs(10)
    if fish < 5: return 1
    return 2 * flip()
```

- a. (6 points) What are the two smallest numbers that can be returned by a call to flop, and what are the probabilities of each being returned?

Answer. Brute force analysis tells us that the smallest possible return value is 1, and that occurs with probability $P(X < 5) = 10e^{-10.5}$. The second smallest possible return value is twice 3, which happens with probability of $\frac{1-10e^{-10.5}}{3}$.

(The code block for this problem is included here for your convenience)

```
from scipy import stats
from random import choice

def flip():
    eel = choice([1, 2, 3])
    if eel == 1: return 3
    if eel == 2: return 5 + flip()
    return 11 + flop()

def flop():
    fish = stats.poisson.rvs(10)
    if fish < 5: return 1
    return 2 * flip()
```

- b. (10 points) What are the expected return values of each of the two functions, flip and flop? Please provide an analytical derivation (i.e., do not just run code).

Answer. Some solution here

That's the end of this quiz!