

MTH 201: Probability and Statistics

Section A End Semester Exam

13/06/2023

Sanjit K. Kaul

No books, notes, or devices are allowed. Just a pen/ pencil and eraser. Your own work alone must be in your sheet. Explain your answers. Show your steps. Don't waste time on algebra unless straightforward. You have about 120 minutes.

Question 1. 25 marks A teaching assistant (TA) collects a stack of sheets for evaluation. The stack has a total of k sheets. The TA doesn't know the total number of sheets in the stack and must count the sheets in the stack. The TA proceeds sheet by sheet from top of the stack to its bottom. The TA accidentally skips counting any sheet in the stack with probability p independently of other sheets. Answer the following questions.

- (10 marks) The TA reports a count of the number of sheets after counting the stack. Derive the PMF of the count.
- (5 marks) Suppose the TA repeats the process of counting sheets in the stack, one or more times, till the TA gets the correct count. Derive the PMF of the total number of times the TA carries out the process of counting sheets in the stack.
- (10 marks) Suppose the TA decides to carry out the process of counting the number of sheets in the stack a total of $m > 1$ times. Having obtained m counts, the TA chooses the maximum of the m counts to be the correct count of sheets in the stack. Derive the PMF of this maximum.

Question 2. 25 marks Lazy must spend 8 hours at work. He decides to spend 3 of the 8 hours on social media (SM) and 5 hours waiting for the weekend (WW). At the beginning of any hour, Lazy chooses either SM or WW randomly from the number of SM and WW hours remaining in the day. Let us index the hours at work as $1, 2, \dots, 8$. The RV $X_i = 1$, $i \in \{1, 2, \dots, 8\}$, if Lazy chooses WW in the i^{th} hour. Else, $X_i = 0$. Answer the following questions.

- (5 marks) Derive the marginal PMFs of the RVs X_i , $i = 1, 2, \dots, 8$. [Hint: Think counting/ combinatorics]
- (5 marks) Derive the joint PMFs of the RVs X_i and X_j , for $i \neq j$.
- (2.5 marks) Are the RVs X_1, X_2, \dots, X_8 mutually independent? Justify your answer.
- (2.5 marks) Derive the expected value of the sum $X_1 + X_2 + \dots + X_8$.
- (5 marks) Derive the variance of the sum $X_1 + X_2 + \dots + X_8$.
- (5 marks) Suppose Lazy chooses to spend at least one hour on SM before spending the first hour on WW. Calculate the conditional expected value of X_3 .

Question 3. 20 marks A virus is known to mutate into variants over months. Any person is infected by the virus within 0 or more integer number of months of the virus having been first detected. The number of months is distributed as a Poisson random variable with expected value $\alpha = 1$ and PMF

$\alpha^k e^{-\alpha}/k!$, $k = 0, 1, \dots$. A person infected by the virus within m months, $m = 0, 1, \dots$, takes an additional exponentially distributed months with rate $\lambda_m = m + 1$ to show symptoms of infection. The PDF of the exponential RV is $\lambda_m e^{-\lambda_m x}$, $x \geq 0$, and its expected value is $1/\lambda_m$. Derive the expected value of the months, since the virus is first detected, any person takes to show symptoms of infection by the virus.

Question 4. 10 marks You arrive at a bus stop knowing that the time you must wait for a bus is distributed as an exponential RV with PDF $f_X(x) = \lambda e^{-\lambda x}$, $x \geq 0$. The bus hasn't arrived for the first s seconds of your wait. Derive the conditional CDF of X . Repeat the above under the assumption that you must wait for a bus for a time that is uniformly distributed over $(0, 120)$. Note that $s \in (0, 120)$.

Question 5. 20 marks An external agency is visiting IIIT-Delhi to evaluate students at the university. A student's evaluation results in a score that is a sum of the marks obtained by the student in humanities and the marks obtained in sciences. The agency would like the average score, where the average is calculated over scores of all students in the institute. Given the large number of students, the agency decides to instead estimate the average. The agency is able to conduct the sciences exam for 50 students and the humanities exam for 25 students. They ask IIIT-Delhi to provide lists of randomly chosen students that will take the two exams.

Assume that the marks of any student i in the sciences exam is given by RV X_i and marks in the humanities exam is given by Y_i . Marks obtained by students in sciences are independent of other marks and identically distributed as the RV X . Marks obtained by students in humanities are independent of other marks and identically distributed as the RV Y . Propose an unbiased estimator for the average. You must show that your estimator is unbiased.

Assume that $\text{Var}[X] = \text{Var}[Y] = 25$. Apply the Chebyshev's inequality to derive the confidence interval $2c$ for a confidence of 95%. The Chebyshev's inequality states that for any RV Z , $P[|Z - E[Z]| \geq c] \leq \text{Var}[Z]/c^2$, for any $c > 0$.

Question 1. 25 marks A teaching assistant (TA) collects a stack of sheets for evaluation. The stack has total of k sheets. The TA doesn't know the total number of sheets in the stack and must count the sheets in the stack. The TA proceeds sheet by sheet from top of the stack to its bottom. The TA accidentally skips counting any sheet in the stack with probability p independently of other sheets. Answer the following questions.

(a) (10 marks) The TA reports a count of the number of sheets after counting the stack. Derive the PMF of the count.

(b) (5 marks) Suppose the TA repeats the process of counting sheets in the stack, one or more times, till the TA gets the correct count. Derive the PMF of the total number of times the TA carries out the process of counting sheets in the stack.

(c) (10 marks) Suppose the TA decides to carry out the process of counting the number of sheets in the stack a total of $m > 1$ times. Having obtained m counts, the TA chooses the maximum of the m counts to be the correct count of sheets in the stack. Derive the PMF of this maximum.

(a) The TA skips a sheet with probability p independently of other sheets. The total no. of sheets is k .

The TA will come up with a count in the set $\{0,1,2,...,k\}$

Let the count be X .
 $S_X = \{0,1,2,...,k\}$

$$P[X=0] = P[\text{TA skips all } k \text{ sheets}] = p^k.$$

$$P[X=1] = \binom{k}{1} p^{k-1} (1-p)$$

\vdots

$$\text{We have } P[X=x] = \begin{cases} \binom{k}{x} p^x (1-p)^{k-x} & x \in \{0,1,2,...,k\} \\ 0 & \text{otherwise.} \end{cases}$$

As you may have guessed, X is a Binomial($k, (1-p)$) RV.

(b) The TA must count the stack again & again till the first time the count is k .

The success event of interest occurs when the TA doesn't skip any sheet when counting. The probability of the event is $(1-p)^k$.

Let Y be the no. of times the TA must count the stack. Let $q = (1-p)^k$.

$$S_Y = \{1,2,...\}$$

$$P[Y=1] = q$$

$$P[Y=2] = (1-q)q$$

$$P[Y=y] = \begin{cases} (1-q)^{y-1} q & y=1,2,... \\ 0 & \text{otherwise.} \end{cases}$$

$$= \begin{cases} (1-(1-p)^k)^{y-1} (1-p)^k & y=1,2,... \\ 0 & \text{otherwise.} \end{cases}$$

Y is geometric($(1-p)^k$).

(c) Let $Z_1, Z_2, ..., Z_m$ be the m counts. Each is a Binomial($k, (1-p)$) RV. We want the PMF of

$$Z = \max(Z_1, Z_2, ..., Z_m)$$

Note that Z has a range space $S_Z = \{0,1,...,k\}$.

$$\begin{aligned} P[Z=0] &= P[\max(Z_1, Z_2, ..., Z_m) = 0] \\ &= P[Z_1=0, Z_2=0, ..., Z_m=0] \\ &= (P[Z_1=0]) (P[Z_2=0]) \dots (P[Z_m=0]) \quad \left[\begin{array}{l} \because Z_i \\ \text{are indep. RVs} \end{array} \right] \\ &= (p^k) (p^k) \dots (p^k) \\ &= (p^k)^m \end{aligned}$$

$$\begin{aligned} P[Z=1] &= P[\max(Z_1, Z_2, ..., Z_m) = 1] \\ &= P[(m-1) \text{ of the } Z_i\text{'s are } \leq 1 \text{ \& one is exactly } 1] \\ &= m \binom{m-1}{1} (P[Z_i=1])^{m-1} P[Z_i=1] \end{aligned}$$

$$P[Z=z] = \begin{cases} m \binom{m-1}{z} (P[Z_i=z])^{m-1} P[Z_i=z] & z=1,2,...,k \\ (P[Z_i=0])^m = (p^k)^m & z=0 \\ 0 & \text{otherwise} \end{cases}$$

At most ①

Some examples At most ②

At most ⑤

At most ①.

At most ①

Some examples At most ①

At most ⑤

At most ①

At most

②

①

Some examples of prob calculations (if provided) ③

④

①

At most ⑧

Question 2: 25 marks Lazy must spend 8 hours at work. He decides to spend 3 of the 8 hours on social media (SM) and 5 hours waiting for the weekend (WW). At the beginning of any hour, Lazy chooses either SM or WW randomly from the number of SM and WW hours remaining in the day. Let us index the hours at work as $1, 2, \dots, 8$. The RV $X_i = 1, i \in \{1, 2, \dots, 8\}$, if Lazy chooses WW in the i th hour. Else, $X_i = 0$. Answer the following questions.

(a) (5 marks) Derive the marginal PMFs of the RVs $X_i, i = 1, 2, \dots, 8$. [Hint: Think counting/ combinatorics]

(b) (5 marks) Derive the joint PMFs of the RVs X_i and X_j , for $i \neq j$.

(c) (2.5 marks) Are the RVs X_1, X_2, \dots, X_8 mutually independent? Justify your answer.

(d) (2.5 marks) Derive the expected value of the sum $X_1 + X_2 + \dots + X_8$.

(e) (5 marks) Derive the variance of the sum $X_1 + X_2 + \dots + X_8$.

(f) (5 marks) Suppose Lazy chooses to spend at least one hour on SM before spending the first hour on WW. Calculated the conditional expected value of X_3 .

(a)
$$P[X_i=1] = P[\text{Lazy chooses WW in the } i\text{th hour and the other hours have a total of 4 WW and 3 SM}]$$

No. of ways in which lazy can assign WW on SM to the hours with the i th hour fixed for WW is $7C_3$.

Total no. of ways of assigning WW on SM to the hours is $8C_3$.

$$\therefore P[X_i=1] = \frac{7C_3}{8C_3} = \frac{\binom{7}{3}}{\binom{8}{3}} = \frac{35}{56} = \frac{5}{8} //$$

Thus the PMF of X_i is

$$P[X_i=x] = \begin{cases} 3/8 & x=0 \\ 5/8 & x=1 \\ 0 & \text{otherwise} \end{cases}$$

Note that this is true for all i in $\{1, 2, \dots, 8\}$ since the prob is independent of i .

(b) Consider

$$P[X_i=1, X_j=1]$$

= P[the i th & j th hour are spent on WW and the remaining 3 WW & 3 SM]

$$= \frac{6C_3}{8C_3} = \frac{6}{8} \cdot \frac{35}{56} = \frac{6}{8} \cdot \frac{5}{8} = \left(\frac{5}{8}\right)\left(\frac{5}{8}\right) = \left(\frac{5}{8}\right)^2$$

$$P[X_i=0, X_j=1] = \frac{6C_2}{8C_3} = \left(\frac{3}{8}\right)\left(\frac{5}{8}\right)$$

$$P[X_i=1, X_j=0] = \frac{6C_2}{8C_3} = \left(\frac{3}{8}\right)\left(\frac{5}{8}\right)$$

$$P[X_i=0, X_j=0] = \frac{6C_1}{8C_3} = \left(\frac{3}{8}\right)\left(\frac{3}{8}\right)$$

The PMF contributes the above four probabilities. Of course for any other pair of values, the probability is 0.

(c) It is easy to show that the RV(s) are not independent. Pick any one of the joint probabilities in (b).

Say $P[X_i=0, X_j=0] = \left(\frac{3}{8}\right)\left(\frac{3}{8}\right)$

We know from (a) that $P[X_i=0] P[X_j=0] = \left(\frac{3}{8}\right)\left(\frac{3}{8}\right)$

Thus, the RV(s) are not independent

(d) $E[X_1 + X_2 + \dots + X_8]$

$$= E[X_1] + E[X_2] + \dots + E[X_8]$$

$$E[X_i] = (1)\left(\frac{5}{8}\right) + (0)\left(\frac{3}{8}\right) = \frac{5}{8}$$

$$\therefore E[X_1 + \dots + X_8] = 5$$

(e) $\text{Var}[X_1 + \dots + X_8]$

$$= \text{Sum of elements of the covariance matrix.}$$

$$= \text{Var}[X_1] + \dots + \text{Var}[X_8]$$

$$+ \text{Cov}[X_1, X_2] + \dots + \text{Cov}[X_1, X_8]$$

$$+ \text{Cov}[X_8, X_1] + \dots + \text{Cov}[X_8, X_7]$$

$$= 8 \text{Var}[X_1] + (8^2 - 8) \text{Cov}[X_1, X_2]$$

$$\text{Var}[X_i] = E[X_i^2] - (E[X_i])^2 = (1)^2 \left(\frac{5}{8}\right) - \left(\frac{5}{8}\right)^2 = \frac{5}{8} \left(\frac{3}{8}\right)$$

$$\text{Cov}[X_1, X_2] = E[X_1, X_2] - E[X_1]E[X_2]$$

$$E[X_1, X_2] = (1)(1)P[X_1=1, X_2=1] = (1)(1)\left(\frac{5}{8}\right)\left(\frac{4}{8}\right) = \frac{20}{64}$$

$$E[X_1]E[X_2] = \left(\frac{5}{8}\right)\left(\frac{5}{8}\right) = \frac{25}{64}$$

$$\therefore \text{Var}[X_1 + \dots + X_8] = 8 \left(\frac{5}{8}\right)\left(\frac{3}{8}\right) + 8(8-1)\left[\frac{20}{64} - \frac{25}{64}\right] = \frac{15}{8} + \left[20 - \frac{35 \times 7}{8}\right] = \frac{15}{8} + \left[\frac{160 - 175}{8}\right] = 0$$

(f) Let A be the event that: lazy chooses to spend at least one hour on SM before spending an hour on WW.

$$E[X_3|A] = (1)P[X_3=1|A] + (0)P[X_3=0|A] = P[X_3=1|A]$$

$$\frac{P[X_3=1, A]}{P[A]} = \frac{P[X_3=1, X_1=0, X_2=0] + P[X_3=1, X_1=0, X_2=1]}{P[A]}$$

$$= \frac{\left(\frac{3}{8}\right)\left(\frac{3}{8}\right)\left(\frac{5}{8}\right) + \left(\frac{3}{8}\right)\left(\frac{5}{8}\right)\left(\frac{4}{8}\right)}{\left[\left(\frac{3}{8}\right)\left(\frac{5}{8}\right) + \left(\frac{3}{8}\right)\left(\frac{2}{8}\right)\left(\frac{5}{8}\right) + \left(\frac{3}{8}\right)\left(\frac{2}{8}\right)\left(\frac{1}{8}\right)\right]}$$

Number Numerator: Of course alternate methods may be used.

Valid arguments for why $5/8$ is the prob are okay too.

At least 1 for any attempt

1 If in line with what was done in (b)

1 If in line with what was done in (a)

0.5 Final result

At least 1 for the expansion

2.5 if done right. Wrong prob calc doesn't impact marks for this part.

Using intuition & getting the correct answer is fine too. Note that the sum $X_1 + \dots + X_8 = 5$ w.p. 1.

Uphs 2

Note that this simplification is bc cause PMF of X_i is the same for all i so one the joint PMF(s) of the pairs

Uphs 1

Uphs

1

Using intuition & getting the correct answer is fine too. Since the sum $X_1 + \dots + X_8 = 5$ always, $\text{Var}[\sum] = 0$.

Expansion of Numerator 1

Denom 1

Correctness 1

Other ways are likely.

Question 3. 20 marks A virus is known to mutate into variants over months. Any person is infected by the virus within 0 or more integer number of months of the virus having been first detected. The number of months is distributed as a Poisson random variable with expected value $\alpha = 1$ and PMF $\alpha^k e^{-\alpha} / k!$, $k = 0, 1, \dots$. A person infected by the virus within m months, $m = 0, 1, \dots$, takes an additional exponentially distributed months with rate $\lambda_m = m + 1$ to show symptoms of infection. The PDF of the exponential RV is $\lambda_m e^{-\lambda_m x}$, $x \geq 0$, and its expected value is $1/\lambda_m$. Derive the expected value of the months, since the virus is first detected, any person takes to show symptoms of infection by the virus.

let X be the total months for a person to show symptoms

let M be the no. of months within which any person is infected.

$$E[X] = \sum_{m=0}^{\infty} E[X | M=m] P[M=m]$$

Up for 5

$$E[X | M=m] = m + \frac{1}{\lambda_m} = m + \frac{1}{m+1}$$

This is for infection

Average for showing symptoms once infected

The conditional

Up for 10

Similar split in case people go the route of calculating the PDF of X .

$$E[X] = \sum_{m=0}^{\infty} \left(m + \frac{1}{m+1} \right) \frac{e^{-1}}{m!}$$

$$= \sum_{m=0}^{\infty} m \frac{e^{-1}}{m!} + \sum_{m=0}^{\infty} \frac{e^{-1}}{(m+1)!}$$

Expected value of Poisson ($\alpha=1$)

$$= 1 + e^{-1} \left[1 + \frac{1}{2} + \frac{1}{2} + \dots \right]$$

$$= 1 + e^{-1} [e - 1]$$

$$= 1 + \left(1 - \frac{1}{e} \right) = 2 - \frac{1}{e}$$

Up for 5

Question 4. 10 marks You arrive at a bus stop knowing that the time you must wait for a bus is distributed as an exponential RV with PDF $f_X(x) = \lambda e^{-\lambda x}, x \geq 0$. The bus hasn't arrived for the first s seconds of your wait. Derive the conditional CDF of X . Repeat the above under the assumption that you must wait for a bus for a time that is uniformly distributed over $(0, 120)$. Note that $s \in (0, 120)$.

We want

$$P[X \leq x | X > s]$$

$$= 1 - P[X > x | X > s]$$

$$P[X > x | X > s] = \begin{cases} \frac{P[X > x]}{P[X > s]} & x > s \\ 1 & \text{otherwise} \end{cases}$$

Up to
(2)

When $f_X(x)$ is exponential,

$$P[X > x | X > s] = \begin{cases} \frac{e^{-\lambda x}}{e^{-\lambda s}} & x > s \\ 1 & \text{otherwise} \end{cases}$$

Up to
(4)

$$\therefore P[X \leq x | X > s] = \begin{cases} 1 - e^{-\lambda(x-s)} & x > s \\ 0 & \text{otherwise} \end{cases}$$

When $f_X(x)$ is uniform $(0, 120)$

$$P[X > x | X > s] = \begin{cases} \frac{120-x}{120-s} & x > s \\ 1 & \text{otherwise} \end{cases}$$

Up to
(4)

$$\Rightarrow P[X \leq x | X > s] = \begin{cases} 1 - \frac{(120-x)}{(120-s)} & x > s \\ 0 & \text{otherwise} \end{cases}$$

Question 5. 20 marks An external agency is visiting IIT-Delhi to evaluate students at the university. A student's evaluation results in a score that is a sum of the marks obtained by the student in humanities and the marks obtained in sciences. The agency would like the average score, where the average is calculated over scores of all students in the institute. Given the large number of students, the agency decides to instead estimate the average. The agency is able to conduct the sciences exam for 50 students and the humanities exam for 25 students. They ask IIT-Delhi to provide lists of randomly chosen students that will take the two exams.

Assume that the marks of any student i in the sciences exam is given by RV X_i and marks in the humanities exam is given by Y_i . Marks obtained by students in sciences are independent of other marks and identically distributed as the RV X . Marks obtained by students in humanities are independent of other marks and identically distributed as the RV Y . Propose an unbiased estimator for the average. You must show that your estimator is unbiased.

Assume that $\text{Var}[X] = \text{Var}[Y] = 25$. Apply the Chebyshev's inequality to derive the confidence interval $2c$ for a confidence of 95%. The Chebyshev's inequality states that for any RV Z , $P[|Z - E[Z]| \geq c] \leq \text{Var}[Z]/c^2$, for any $c > 0$.

We want to estimate $E[X + Y]$.
The unbiased estimator is

$$\begin{aligned} Z &= M_{n_1}(X) + M_{n_2}(Y) \\ &= \frac{1}{n_1} \sum_{i=1}^{n_1} X_i + \frac{1}{n_2} \sum_{i=1}^{n_2} Y_i \end{aligned} \quad \left\{ \begin{array}{l} n_1=50 \\ n_2=25 \end{array} \right.$$

This can be seen by taking its expectation.

$$\begin{aligned} E[M_{n_1}(X) + M_{n_2}(Y)] &= \frac{1}{n_1} \sum_{i=1}^{n_1} E[X] + \frac{1}{n_2} \sum_{i=1}^{n_2} E[Y] \\ &= E[X] + E[Y]. \end{aligned}$$

$$P[|Z - E[Z]| \geq c] \leq \frac{\text{Var}[Z]}{c^2}$$

$$\begin{aligned} \text{Var}[Z] &= \text{Var}[M_{n_1}(X) + M_{n_2}(Y)] \\ &= \text{Var}[M_{n_1}(X)] + \text{Var}[M_{n_2}(Y)] \\ &= \frac{1}{n_1} \text{Var}[X] + \frac{1}{n_2} \text{Var}[Y] \\ &= \frac{25}{50} + \frac{25}{25} \end{aligned}$$

We want 95% confidence.

That is choose c such that

$$\frac{\text{Var}[Z]}{c^2} \leq 0.05.$$

$$\begin{aligned} c^2 &= \frac{\text{Var}[Z] (100)}{5} = \left(\frac{1}{2} + 1\right) (20) \\ &= 30 \end{aligned}$$

satisfies the requirement.

$$c = \sqrt{30}.$$