

# Spring 2025 CS4641/CS7641 Homework 1

Dong Kyun Lim

Deadline: Friday, February 7th, 11:59 pm EST

- No unapproved extension of the deadline is allowed. For late submissions, please refer to the course website.
- Discussion is encouraged on Ed as part of the Q/A. However, all assignments should be done individually.
- Plagiarism is a **serious offense**. You are responsible for completing your own work. You are not allowed to copy and paste, or paraphrase, or submit materials created or published by others, as if you created the materials. All materials submitted must be your own. This also means you may not submit work created by generative models as your own.
- All incidents of suspected dishonesty, plagiarism, or violations of the Georgia Tech Honor Code will be subject to the institute's Academic Integrity procedures. If we observe any (even small) similarities/plagiarisms detected by Gradescope or our TAs, **WE WILL DIRECTLY REPORT ALL CASES TO OSI**, which may, unfortunately, lead to a very harsh outcome. **Consequences can be severe, e.g., academic probation or dismissal, grade penalties, a 0 grade for assignments concerned, and prohibition from withdrawing from the class.**

## Instructions

- We will be using Gradescope for submission and grading of assignments.
- Unless a question explicitly states that no work is required to be shown, you must provide an explanation, justification, or calculation for your answer. Basic arithmetic can be combined (it does not need to each have its own step); your work should be at a level of detail that a TA can follow it.
- Your write-up must be submitted in PDF form, you may use either Latex, markdown, or any word processing software. **We will NOT accept handwritten work.** Make sure that your work is formatted correctly, for example submit  $\sum_{i=0} x_i$  instead of `sum_{i=0} x.i`.
- A useful video tutorial on LaTeX has been created by our TA team and can be found [here](#) and an Overleaf document with the commands can be found [here](#).
- When submitting your assignment on Gradescope, you are required to correctly map pages of your PDF to each question/ subquestion to reflect where they appear. Improperly mapped questions will not be graded correctly.
- All assignments should be done individually, each student must write up and submit their own answers.
- **Graduate Students:** You are required to complete any sections marked as Bonus for Undergrads

## Point Distribution

### Q1: Linear Algebra [30pts]

- 1.1 Determinant and Inverse of a Matrix [10pts]
- 1.2 Eigenvalues and Eigenvectors [20pts]

### Q2: Expectation, Co-variance and Statistical Independence [7pts]

### Q3: Optimization [17pts + 3% Bonus for All]

- 3.1 KKT [17pts]
- 3.2 Primal and Dual Form [3% Bonus for All]

### Q4: Maximum Likelihood [20pts: 10pts + 10 pts Grad/6% Bonus for Undergrads]

- 4.1 Discrete Example [10pts]
- 4.2 Poisson Distribution [10pts Grad / 6% Bonus for Undergrads]

### Q5: Information Theory [31pts]

- 5.1 Mutual Information and Entropy [21pts]
- 5.2 Entropy Proofs [10pts]

### Q6: Ethical Implications on Decision-Making [10 pts]

- 6.1 Loan Eligibility [5pts]
- 6.2 Voting and Probabilistic Models [5pts]

### Q7: Programming [5pts]

### Q8: Bonus Questions [7% Bonus for All]

- 8.1 Marginal Probability Density Functions [2% Bonus for All]
- 8.2 Coin Toss Game [2% Bonus for All]
- 8.3 Dice Roll Expectation [3% Bonus for All]

### Points Totals:

- **Total Programming Points for All:** 5 pts
- **Total Written Points for Grad:** 115 pts
- **Total Written Points for Undergrad:** 105 pts

# 1 Linear Algebra [30pts]

## 1.1 Determinant and Inverse of Matrix [10pts]

Given a matrix  $M$ :

$$M = \begin{bmatrix} 3 & -2 & 4 \\ r & 1 & -1 \\ 0 & 2 & 2 \end{bmatrix}$$

- (a) Calculate the determinant of  $M$  in terms of  $r$  (calculation process is required). [4pts]

$$|M| = 3 \cdot \begin{vmatrix} 1 & -1 \\ 2 & 2 \end{vmatrix} - (-2) \cdot \begin{vmatrix} r & -1 \\ 0 & 2 \end{vmatrix} + 4 \cdot \begin{vmatrix} r & 1 \\ 0 & 2 \end{vmatrix}$$

$$3 \cdot ((1)(2) - (-1)(2)) = 12$$

$$-2 \cdot ((r)(2) - (-1)(0)) = -4r$$

$$4 \cdot ((r)(2) - (1)(0)) = 8r$$

$$12 - (-4r) + 8r = 12 + 12r$$

$$\therefore |M| = 12 + 12r$$

- (b) For what value(s) of  $r$  does  $M^{-1}$  not exist? Why doesn't  $M^{-1}$  exist in this case? What does it mean in terms of rank and singularity for these values of  $r$ ? *This question can be answered in less than 7 lines.* [3pts]

For the inverse of  $M$  to exist, the determinant must be non-zero. Therefore, when the determinant is zero, the inverse does not exist.

$$|M| = 12 + 12r = 0 \Rightarrow r = -1$$

If the determinant is zero, the matrix is called **singular** which makes the columns of the matrix linearly dependent meaning at least one column can be expressed as a linear combination of others. This would also mean that the rank of the matrix is strictly less than 3 (full rank).

- (c) Find the mathematical equation that describes the relationship between the determinant of  $M$  and the determinant of  $M^{-1}$ . You do **NOT** need to show any work. [3pts]

**NOTE:** It may be helpful to find the determinant of  $M$  and  $M^{-1}$  for  $r = 0$ .

$$|M^{-1}| = \frac{1}{|M|}$$

## 1.2 Eigenvalues and Eigenvectors [20pts]

### 1.2.1 Eigenvalues [5pts]

Given the following matrix  $\mathbf{A}$ , find an expression for the eigenvalues  $\lambda$  of  $\mathbf{A}$  in terms of  $a$ ,  $b$ , and  $c$ . (Simplify your answer into the form  $\lambda = \dots$ ). [5pts]

$$\mathbf{A} = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

$$\mathbf{A} - \lambda \mathbf{I} = \begin{bmatrix} a & b \\ b & c \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = \begin{bmatrix} a - \lambda & b \\ b & c - \lambda \end{bmatrix}$$

$$\begin{vmatrix} a - \lambda & b \\ b & c - \lambda \end{vmatrix} = (a - \lambda)(c - \lambda) - b^2 = \lambda^2 - (a + c)\lambda + ac - b^2$$

$$\lambda^2 - (a + c)\lambda + ac - b^2 = 0$$

Use quadratic formula:

$$\therefore \lambda = \frac{a + c \pm \sqrt{(a + c)^2 - 4(ac - b^2)}}{2}$$

### 1.2.2 Eigenvectors [15pts]

Given a matrix  $\mathbf{A}$ :

$$\mathbf{A} = \begin{bmatrix} -7 & 2 \\ 6 & 4 \end{bmatrix}$$

(a) Calculate the eigenvalues of  $\mathbf{A}$ . Simplify your answer into the form  $\lambda = \text{numbers}$  [3pts]

$$\mathbf{A} - \lambda \mathbf{I} = \begin{bmatrix} -7 & 2 \\ 6 & 4 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = \begin{bmatrix} -7 - \lambda & 2 \\ 6 & 4 - \lambda \end{bmatrix}$$

$$\begin{vmatrix} -7 - \lambda & 2 \\ 6 & 4 - \lambda \end{vmatrix} = (-7 - \lambda)(4 - \lambda) - 12 = \lambda^2 + 3\lambda - 40 = 0$$

$$(\lambda - 5)(\lambda + 8) = 0$$

$$\therefore \lambda = -8, 5$$

(b) Find the normalized eigenvectors of matrix  $\mathbf{A}$  (calculation process required). [7pts]

$$\mathbf{A} - 5\mathbf{I} = \begin{bmatrix} -7 - 5 & 2 \\ 6 & 4 - 5 \end{bmatrix} = \begin{bmatrix} -12 & 2 \\ 6 & -1 \end{bmatrix}$$

$$\begin{bmatrix} -12 & 2 \\ 6 & -1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$-12x + 2y = 0$$

$$y = 6x$$

$$\|v_1\| = \sqrt{1^2 + 6^2} = \sqrt{37}$$

$$\therefore v_1 = \frac{1}{\sqrt{37}} \begin{bmatrix} 1 \\ 6 \end{bmatrix}$$

$$\mathbf{A} + 8\mathbf{I} = \begin{bmatrix} -7 + 8 & 2 \\ 6 & 4 + 8 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 6 & 12 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 \\ 6 & 12 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$x + 2y = 0$$

$$x = -2y$$

$$||v_2|| = \sqrt{1^2 + 2^2} = \sqrt{5}$$

$$\therefore v_2 = \frac{1}{\sqrt{5}} \begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

- (c) When calculating the eigenvectors, were the columns of the matrix  $(\mathbf{A} - \lambda \mathbf{I})$  linearly independent or linearly dependent?

The columns were linearly dependent.

Now, consider the linearly independent matrix

$$\mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$$

Solve for the vector  $\mathbf{x}$  which satisfies the equation  $\mathbf{B}\mathbf{x} = \mathbf{0}$ .

Hint: Use row reduction on

$$\left[ \begin{array}{cc|c} 1 & 2 & 0 \\ 2 & 1 & 0 \end{array} \right]$$

$$R_2 = R_2 - 2R_1 \Rightarrow \begin{bmatrix} 1 & 2 \\ 0 & -3 \end{bmatrix}$$

$$R_2 = \frac{-R_2}{3} \Rightarrow \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix}$$

$$R_1 = R_1 - 2R_2 \Rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\therefore x_1 = 0, x_2 = 0$$

Afterwards, recall that matrices can be interpreted as a transformation on a vector. For example,

$$\text{scaling} = \begin{bmatrix} k & 0 \\ 0 & k \end{bmatrix}, \text{flip across } x_2 = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$$

So, consider you have several vectors close to each other, and you want to apply a transformation to separate them. If you want to make sure the vectors keep their non-zero values after the transformation, what important property must the transformation matrix have [5pts]?

The transformation matrix must be invertible to prevent vectors from collapsing to zero.

## 2 Expectation, Co-variance, and Statistical Independence [7pts]

Suppose  $X$ ,  $Y$ , and  $Z$  are three different real-valued random variables.

Let  $X$  obey a discrete binary distribution. The probability mass function for  $X$  is:

$$p(x) = \begin{cases} 0.8 & x = c \\ 0.2 & x = -c \end{cases}$$

where  $c$  is some nonzero constant. The distribution of  $Y$  is not known, but it is provided that  $\text{Var}(Y) = 0.94c^2$ . Additionally,  $X$  and  $Y$  are statistically independent (i.e.  $P(X|Y) = P(X)$ ). Finally, let  $Z = 9X + 3Y$ .

We define a correlational measure  $\gamma$ :

$$\gamma(X, Z) = \frac{\text{Cov}(X, Z)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Z)}} + \text{Var}(X + Z)$$

Evaluate  $\gamma(X, Z)$  in terms of  $c$ . Remember to show your work to receive credit. Round the values in your final answer to 3 decimal places, but do not round in intermediate steps.

**HINT:** Review the probability and statistics lecture slides for relevant formulae.

$$E[X] = c \cdot 0.8 - c \cdot 0.2 = 0.6c$$

$$E[X^2] = c^2 \cdot 0.8 + c^2 \cdot 0.2 = c^2$$

$$\text{Var}(X) = E[X^2] - (E[X])^2 = c^2 - (0.6c)^2 = 0.64c^2$$

$$\text{Cov}(X, Z) = \text{Cov}(X, 9X + 3Y) = 9\text{Cov}(X, X) + 3\text{Cov}(X, Y)$$

Since,  $X$  and  $Y$  are independent, covariance between  $X$  and  $Y$  is 0.

$$\text{Cov}(X, Z) = 9\text{Cov}(X, X) = 9\text{Var}(X) = 9(0.64c^2) = 5.76c^2$$

$$\text{Var}(Z) = \text{Var}(9X + 3Y) = 9^2\text{Var}(X) + 3^2\text{Var}(Y) = 81(0.64c^2) + 9(0.94c^2) = 60.3c^2$$

$$\text{Var}(X + Z) = \text{Var}(X) + \text{Var}(Z) + 2\text{Cov}(X, Z)$$

$$\text{Var}(X + Z) = 0.64c^2 + 60.3c^2 + 2(5.76c^2) = 72.46c^2$$

$$\therefore \gamma(X, Z) = \frac{5.76c^2}{\sqrt{(0.64c^2)(60.3c^2)}} + 72.46c^2 = \frac{5.76}{\sqrt{38.592}} + 72.46c^2 = 0.927 + 72.46c^2$$

### 3 Optimization [17pts + 3% Bonus for All]

#### 3.1 KKT [17pts]

Optimization problems are related to minimizing a function (usually termed loss, cost or error function) or maximizing a function (such as the likelihood) with respect to some variable  $x$ . The Karush-Kuhn-Tucker (KKT) conditions are first-order conditions for a solution in nonlinear programming to be optimal, provided that some regularity conditions are satisfied.

In this question, you will be solving the following optimization problem. In this problem, you are tasked with helping professor Mahdi optimise the type and number of GPU's to buy. You must balance cost and availability. You have the total compute defined by function  $f(x,y)$ , and the constraints with respect to cost and availability.

$$\begin{aligned} \max_{x,y} \quad & f(x,y) = 7x^2 + 4y^2 \\ \text{s.t.} \quad & g_1(x,y) = 15x + 6y \leq 250 \\ & g_2(x,y) = x \leq 8 \end{aligned}$$

- (a) Write the Lagrange function for the maximization problem. Now change the maximum function to a minimum function (i.e.  $\min_{x,y} f(x,y) = 7x^2 + 4y^2$ ) and provide the Lagrange function for the minimization problem with the same constraints  $g_1$  and  $g_2$ . [2pts]

**NOTE:** The minimization problem is only for part (a).

For maximization:

$$\mathcal{L}(x, y, \lambda_1, \lambda_2) = 7x^2 + 4y^2 - \lambda_1(15x + 6y - 250) - \lambda_2(x - 8)$$

For minimization:

$$\mathcal{L}(x, y, \lambda_1, \lambda_2) = 7x^2 + 4y^2 + \lambda_1(15x + 6y - 250) + \lambda_2(x - 8)$$

- (b) List the names of all 4 groups of KKT conditions and their corresponding mathematical equations or inequalities for this specific maximization problem. Be sure to simplify completely, and calculate the derivative. [2pts]

1. Stationary Conditions: Compute partial derivatives

$$\frac{\partial \mathcal{L}}{\partial x} = 14x - 15\lambda_1 - \lambda_2 = 0$$

$$\frac{\partial \mathcal{L}}{\partial y} = 8y - 6\lambda_1 = 0$$

2. Complementary Slackness

$$\lambda_1(15x + 6y - 250) = 0$$

$$\lambda_2(x - 8) = 0$$

3. Primal Feasibility

$$15x + 6y - 250 \leq 0$$

$$x - 8 \leq 0$$

4. Dual Feasibility

$$\lambda_1, \lambda_2 \geq 0$$

- (c) Solve for 4 possibilities formed by each constraint being active or inactive. Do not forget to check the inactive constraints for each point when applicable. Candidate points must satisfy all the conditions mentioned in part b) (Quick note: If a constraint is binding its corresponding lambda should be greater than or equal to zero). [8pts]

*Case1 :  $g_1 \rightarrow \text{active}, g_2 \rightarrow \text{active}$*

$$\lambda_1 \geq 0$$

$$\lambda_2 \geq 0$$

$$x - 8 = 0 \Rightarrow x = 8$$

$$15(8) + 6y = 250$$

$$6y = 130$$

$$y = \frac{65}{3}$$

$$(x, y) = (8, \frac{65}{3})$$

Now, let's use this solution to solve stationary conditions.

$$8(\frac{65}{3}) - 6\lambda_1 = 0$$

$$6\lambda_1 = \frac{520}{3}$$

$$\lambda_1 = \frac{260}{9} \geq 0$$

$$14(8) - 15(\frac{260}{9}) = \lambda_2$$

$$\lambda_2 = 112 - \frac{1300}{3} \leq 0$$

Since one of the lambda does not meet Dual Feasibility, there is no valid point.

*Case2 :  $g_1 \rightarrow \text{active}, g_2 \rightarrow \text{inactive}$*

$$\lambda_1 \geq 0$$

$$\lambda_2 = 0$$

$$x - 8 < 0$$

$$15x + 6y - 250 = 0$$

$$\Rightarrow 6y = 250 - 15x$$

$$\Rightarrow y = \frac{250 - 15x}{6}$$

$$\frac{\partial \mathcal{L}}{\partial x} = 14x - 15\lambda_1 - \lambda_2 = 0$$

$$\Rightarrow 14x - 15\lambda_1 = 0$$

$$\Rightarrow \lambda_1 = \frac{14x}{15}$$

$$\frac{\partial \mathcal{L}}{\partial y} = 8y - 6\lambda_1 = 0$$

$$\Rightarrow y = \frac{3\lambda_1}{4} = 0.7x = \frac{250 - 15x}{6}$$

$$\Rightarrow x = 13.021, y = 9.114$$

Since x is greater than 8, it does not meet Primal Feasibility, therefore no valid point.



Case3 :  $g_1 \rightarrow \text{inactive}, g_2 \rightarrow \text{active}$

$$\lambda_1 = 0$$

$$\lambda_2 \geq 0$$

$$x - 8 = 0 \Rightarrow x = 8$$

$$14(8) - 15(0) - \lambda_2 = 0 \Rightarrow \lambda_2 = 112$$

$$8y - 6(0) = 0 \Rightarrow y = 0$$

$$(x, y) = (8, 0)$$

The solution meets all 4 conditions and therefore, (8,0) is a valid point.

Case4 :  $g_1 \rightarrow \text{inactive}, g_2 \rightarrow \text{inactive}$

$$\lambda_1 = 0$$

$$\lambda_2 = 0$$

$$14x = 0 \Rightarrow x = 0$$

$$8y = 0 \Rightarrow y = 0$$

$$(x, y) = (0, 0)$$

The solution meets all 4 conditions and therefore, (0,0) is a valid point.

- (d) List the candidate point(s) (there is at least 1) obtained from part c). Please round answers to 3 decimal points and use that answer for calculations in further parts. This part can be completed in one line per candidate point. [2pts]

$$(0, 0), \lambda_1 = 0, \lambda_2 = 0$$

$$(8, 0), \lambda_1 = 0, \lambda_2 = 112$$

- (e) Find the **one** candidate point for which  $f(x, y)$  is largest. Check if  $L(x, y)$  is concave, convex, or neither at this point by using the [Hessian](#) in the [second partial derivative test](#). [3pts]

$$f(x, y) = 7x^2 + 4y^2$$

$$f(0, 0) = 7(0)^2 + 4(0)^2 = 0$$

$$f(8, 0) = 7(8)^2 + 4(0)^2 = 448$$

(8,0) maximizes the function.

Now, use Hessian matrix and second partial derivative test to determine concavity of the function.

$$\mathcal{L}(x, y) = 7x^2 + 4y^2 - 112(x - 8)$$

$$H = \begin{bmatrix} \frac{\partial^2 \mathcal{L}}{\partial x^2} & \frac{\partial^2 \mathcal{L}}{\partial x \partial y} \\ \frac{\partial^2 \mathcal{L}}{\partial y \partial x} & \frac{\partial^2 \mathcal{L}}{\partial y^2} \end{bmatrix}$$

Let's evaluate each term of the Hessian matrix.

$$\frac{\partial \mathcal{L}}{\partial x} = 14x - 112$$

$$\frac{\partial \mathcal{L}}{\partial y} = 8y$$

$$\frac{\partial^2 \mathcal{L}}{\partial x^2} = 14$$

$$\frac{\partial^2 \mathcal{L}}{\partial x \partial y} = 0$$

$$\frac{\partial^2 \mathcal{L}}{\partial y \partial x} = 0$$

$$\frac{\partial^2 \mathcal{L}}{\partial y^2} = 8$$

$$H = \begin{bmatrix} 14 & 0 \\ 0 & 8 \end{bmatrix}$$

$$|H| = (14)(8) - (0)(0) = 112$$

Since the determinant is greater than 0 AND the entries of the matrix are positive, Lagrangian is convex.

**HINT:** Read the Example\_optimization\_problem.pdf in Canvas Files for HW1 to see an example with some explanations.

**HINT:** Watch this [video](#) walking you through how to solve a similar problem.

### 3.2 Primal and Dual Form [3% Bonus for All]

Convex optimization problems involve minimizing a function given a constraint. The Lagrangian function includes a penalty for violating this constraint. A maximum is taken over the penalty because this is the opposite of what we are trying to accomplish which is minimization.

Under certain conditions, which are satisfied in the following problem, maximizing a variable followed by minimizing over another variable is equivalent to minimizing over the latter followed by maximizing over the former. In literature, this is referred to as transforming a problem from the primal form into the dual form.

For the following problem, write out the primal form, switch it to the dual form, and then solve for the pair  $(x, y)$ .

**NOTE:** The following [video](#) does a great job at visualizing this concept. Additionally, for linear constraints, there is no "inactive" state for the constraint.

$$\begin{aligned} \min_{x,y} \quad & f(x, y) = x^2 + y^2 \\ \text{s.t.} \quad & g_1(x, y) = x + y = 4 \end{aligned}$$

$$\mathcal{L}(x, y, \lambda) = x^2 + y^2 + \lambda(x + y - 4)$$

Now, we should minimize the above Lagrangian function

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x} = 2x + \lambda = 0 &\Rightarrow x = \frac{-\lambda}{2} \\ \frac{\partial \mathcal{L}}{\partial y} = 2y + \lambda = 0 &\Rightarrow y = \frac{-\lambda}{2} \end{aligned}$$

Now, convert primal form into dual form.

$$\mathcal{L}(\lambda) = \left(\frac{-\lambda}{2}\right)^2 + \left(\frac{-\lambda}{2}\right)^2 + \lambda(-\lambda - 4) = -\frac{\lambda^2}{2} - 4\lambda$$

Now, take derivative of the dual form to find  $\lambda$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \lambda} = -\lambda - 4 &= 0 \\ \lambda &= -4 \\ \Rightarrow (x, y) &= (2, 2) \end{aligned}$$

Now, substitute above values into the function.

$$f(2, 2) = 2^2 + 2^2 = 8$$

## 4 Maximum Likelihood [20pts: 10pts + 10pts Grad / 6% Bonus for Undergrads]

### 4.1 Discrete Example [10pts]

Mastermind Mahdi decides to give a challenge to his students for their MLE Final. He provides a spinner with 10 sections, each numbered 1 through 10. The students can change the sizes of each section, meaning that they can select the probability the spinner lands on a certain section. Mahdi then proposes that the students will get a 100 on their final if they can spin the spinner 10 times such that it doesn't land on section 1 during the first 9 spins and lands on section 1 on the 10th spin. If the probability of the spinner landing on section 1 is  $\theta$ , what value of  $\theta$  should the students select to most likely ensure they get a 100 on their final? Use your knowledge of Maximum Likelihood Estimation to get a 100 on the final.

**NOTE: You must specify the log-likelihood function and use MLE to solve this problem for full credit.** You may assume that the log-likelihood function is concave for this question

2025-1Spring/Solution/spinner\_10.png

Probability of not landing on 1

$$P(\text{Not } 1) = 1 - \theta$$

Probability of landing on 1

$$P(1) = \theta$$

Likelihood of 9 spins not landing on 1 and the 10th spin landing on 1

$$L(\theta) = (1 - \theta)^9 \cdot \theta$$

Log-likelihood function:

$$\begin{aligned} \log L(\theta) &= \log((1 - \theta)^9 \cdot \theta) \\ \Rightarrow \log L(\theta) &= 9\log(1 - \theta) + \log \theta \end{aligned}$$

Maximize log-likelihood function:

$$\begin{aligned} \frac{d}{d\theta}(9\log(1 - \theta) + \log \theta) &= 0 \\ \Rightarrow \frac{-9}{1 - \theta} + \frac{1}{\theta} &= 0 \\ \Rightarrow \frac{1}{\theta} &= \frac{9}{1 - \theta} \\ \Rightarrow 1 - \theta &= 9\theta \\ \Rightarrow 1 &= 10\theta \\ \therefore \theta &= \frac{1}{10} \end{aligned}$$

## 4.2 Normal distribution [10 pts Grad / 6% Bonus for Undergrads]

The Normal distribution is defined as:

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

- (a) Let  $(X_1, \dots, X_n) \sim \mathcal{N}(\mu, \sigma^2)$ , where  $X_1, \dots, X_n$  are i.i.d. random variables, and let  $x_1, \dots, x_n$  be the observed values of  $X_1, \dots, X_n$ . What is the likelihood of  $(\mu, \sigma^2)$  given this data? Express your answer in product form. [4 pts / 2%]

The observations are independent. Therefore, the likelihood is a product of individual likelihoods.

$$L(\mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$L(\mu, \sigma^2) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

- (b) What are the maximum likelihood estimators (MLEs) for  $\mu$  and  $\sigma^2$  (Hint the MLEs of  $\sigma^2$  is in terms of  $\mu$ ) ? [6 pts / 4%]

To find MLE, take natural log to both sides.

$$\ln L(\mu, \sigma^2) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

Now, take derivative with respect to  $\mu$  on both sides.

$$\frac{\partial \ln L}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)$$

Set the derivative to 0.

$$\sum_{i=1}^n (x_i - \hat{\mu}) = 0$$

$$\therefore \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

Now, take derivative with respect to  $\sigma^2$  on both sides.

$$\frac{\partial \ln L}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2$$

Set the derivative to 0.

$$-\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \hat{\mu})^2 = 0$$

Multiply both sides by  $2\sigma^4$

$$-n\sigma^2 + \sum_{i=1}^n (x_i - \hat{\mu})^2 = 0$$

$$\therefore \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

## 5 Information Theory [31pts]

### 5.1 Mutual Information and Entropy [21pts]

A recent study has shown symptomatic infections are responsible for higher transmission rates. Using the data collected from positively tested patients, we wish to determine which feature(s) have the greatest impact on whether or not someone will present with symptoms. To do this, we will compute the entropies, conditional entropies, and mutual information of select features. Please use base 2 when computing logarithms.

ID	Vaccine Doses ( $X_1$ )	Wears Mask? ( $X_2$ )	Underlying Conditions ( $X_3$ )	Symptomatic ( $Y$ )
1	L	T	F	F
2	M	F	F	T
3	L	F	F	F
4	H	T	F	F
5	L	F	T	T
6	H	F	T	T
7	L	T	T	F
8	M	F	F	T
9	H	T	T	F
10	M	T	F	F

Table 1: Vaccine Doses: {(H) booster, (M) 2 doses, (L) 1 dose, (T) True, (F) False}

- (a) Find entropy  $H(Y)$  to at least 3 decimal places. [3pts]

$$\begin{aligned}
 H(Y) &= - \sum P(Y) \log_2 P(Y) \\
 H(Y) &= -(P(Y = T) \log_2 P(Y = T) + P(Y = F) \log_2 P(Y = F)) \\
 H(Y) &= -(0.4) \log_2 0.4 - (0.6) \log_2 0.6 \approx 0.971
 \end{aligned}$$

- (b) Find the average conditional entropy  $H(Y|X_1)$  and  $H(Y|X_2)$  to at least 3 decimal places. [7pts]

$$\begin{aligned}
 H(Y|X) &= \sum P(X) H(Y|X = x) \\
 H(Y|X_1) &= P(X_1 = L) H(Y|X_1 = L) + P(X_1 = M) H(Y|X_1 = M) + P(X_1 = H) H(Y|X_1 = H) \\
 P(X_1 = L) H(Y|X_1 = L) &= \frac{4}{10} \cdot \left( -\left( \frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} \right) \right) \approx 0.324 \\
 P(X_1 = M) H(Y|X_1 = M) &= \frac{3}{10} \cdot \left( -\left( \frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right) \right) \approx 0.275 \\
 P(X_1 = H) H(Y|X_1 = H) &= \frac{3}{10} \cdot \left( -\left( \frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3} \right) \right) \approx 0.275 \\
 \therefore H(Y|X_1) &= 0.324 + 0.275 + 0.275 \approx 0.874 \\
 H(Y|X_2) &= P(X_2 = T) H(Y|X_2 = T) + P(X_2 = F) H(Y|X_2 = F) \\
 P(X_2 = T) H(Y|X_2 = T) &= \frac{5}{10} \cdot 0 = 0 \\
 P(X_2 = F) H(Y|X_2 = F) &= \frac{5}{10} \cdot \left( -\left( \frac{4}{5} \log_2 \frac{4}{5} + \frac{1}{5} \log_2 \frac{1}{5} \right) \right) \approx 0.361 \\
 \therefore H(Y|X_2) &= 0 + 0.361 \approx 0.361
 \end{aligned}$$

- (c) Find mutual information  $I(X_1, Y)$  and  $I(X_2, Y)$  to at least 3 decimal places and determine which one ( $X_1$  or  $X_2$ ) is more informative. [3pts]

$$I(X, Y) = H(Y) - H(Y|X)$$

$$I(X_1, Y) = H(Y) - H(Y|X_1) = 0.971 - 0.874 \approx 0.097$$

$$I(X_2, Y) = H(Y) - H(Y|X_2) = 0.971 - 0.361 \approx 0.610$$

$$I(X_2, Y) > I(X_1, Y)$$

$\therefore X_2$  is more informative.

- (d) Find joint entropy  $H(Y, X_3)$  to at least 3 decimal places. [3pts]

$$H(Y, X_3) = H(Y) + H(X_3|Y)$$

$$\begin{aligned} &= H(Y) + P(Y = T)H(X_3|Y = T) + P(Y = F)H(X_3|Y = F) \\ &= 0.971 + \frac{4}{10} \cdot \left(-\left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4}\right)\right) + \frac{6}{10} \cdot \left(-\left(\frac{2}{6} \log_2 \frac{2}{6} + \frac{4}{6} \log_2 \frac{4}{6}\right)\right) \approx 1.922 \end{aligned}$$

- (e) Find the conditional entropy  $H(Y|X_1, X_2)$ . [5 pts]

$$H(Y|X_1, X_2) = - \sum P(X_1, X_2, Y) \log_2 P(Y|X_1, X_2)$$

Possible  $(X_1, X_2) : (L, T), (L, F), (M, T), (M, F), (H, T), (H, F)$

$$(L, T) : P(X_1 = L, X_2 = T)H(Y|X_1 = L, X_2 = T) = 0$$

$$(L, F) : P(X_1 = L, X_2 = F)H(Y|X_1 = L, X_2 = F) = \frac{2}{10} \cdot 1 = 0.2$$

$$(M, T) : P(X_1 = M, X_2 = T)H(Y|X_1 = M, X_2 = T) = 0$$

$$(M, F) : P(X_1 = M, X_2 = F)H(Y|X_1 = M, X_2 = F) = 0$$

$$(H, T) : P(X_1 = H, X_2 = T)H(Y|X_1 = H, X_2 = T) = 0$$

$$(H, F) : P(X_1 = H, X_2 = F)H(Y|X_1 = H, X_2 = F) = 0$$

$$\therefore H(Y|X_1, X_2) = 0.2$$

## 5.2 Entropy Proofs [10pts]

Given the mathematical definition of  $H(X)$  and  $H(X|Y)$  below, prove that  $I(X, Y) = 0$  if  $X$  and  $Y$  are statistically independent. (Note: you must provide a mathematical proof and cannot use the visualization shown in class [found here](#). You may use any theorem/ proof from the slides without having to re-prove it). [10pts]

$$H(X) = - \sum_x P(x) \log_2 P(x)$$

$$H(X|Y) = - \sum_{x,y} P(x,y) \log_2 \frac{P(x,y)}{P(y)}$$

**Start from:**  $I(X, Y) = H(X) - H(X|Y)$

$$I(X, Y) = H(X) - H(X|Y)$$

$$H(X|Y) = - \sum_{x,y} P(x,y) \log_2 \frac{P(x,y)}{P(y)}$$

$$= - \sum_{x,y} P(x)P(y) \log_2 \frac{P(x)P(y)}{P(y)}$$

$$= - \sum_{x,y} P(x)P(y) \log_2 P(x)$$

$$= - \sum_x P(x) \log_2 P(x) \sum_y P(y)$$

$$\sum_y P(y) = 1$$

$$H(X|Y) = - \sum_x P(x) \log_2 P(x) = H(X)$$

$$\therefore H(X) - H(X|Y) = 0$$



## 6 Ethical Implications on Decision-Making [10 pts]

### 6.1 Loan Eligibility [5pts]

#### Real-world Implications

Loan eligibility determines who can receive a loan, typically based on financial history and demographics. It is a difficult problem, and often uses algorithms to make loan decisions. Often, this can result in reinforcing inequality and bias [oneil].

Suppose we're using a matrix to represent the attributes of individuals for loan approval. Each attribute (like income, credit score, years of employment, etc.) constitutes a column in our matrix. Here's a hypothetical toy example:

	Annual Income	Debt-to-Income Ratio	Employment History (years)	Credit Score
Candidate 1	50,000	0.2	5	700
Candidate 2	51,000	0.21	5.1	710
Candidate 3	45,000	0.19	4.9	690
Candidate 4	100,000	0.05	10	780

One algorithm used to predict credit score is linear regression, formulated as  $\mathbf{y} = \mathbf{x}\mathbf{A}$ .  $\mathbf{y}$  are the target variables,  $\mathbf{x}$  are the input features, and  $\mathbf{A}$  is a matrix trained with an existing dataset. Training data  $(\mathbf{x}_D, \mathbf{y}_D)$  are taken from the training dataset  $D$ ,  $(\mathbf{x}_D, \mathbf{y}_D) \in D$ . If  $\mathbf{x}_D$  is linearly independent,  $\mathbf{A}$  can be trained by simply inverting  $\mathbf{x}_D$ :

$$\begin{aligned}\mathbf{y}_D &= \mathbf{x}_D \mathbf{A} \\ \mathbf{x}_D^{-1} \mathbf{y}_D &= \mathbf{A}\end{aligned}$$

The original equation can be rewritten as:

$$\begin{aligned}\mathbf{y} &= \mathbf{x}\mathbf{A} \\ &= \mathbf{x}\mathbf{x}_D^{-1}\mathbf{y}_D\end{aligned}$$

Problems arise when the training data is close to linearly dependent. Recall that one way to invert a matrix is  $\mathbf{A}^{-1} = \frac{1}{\det(\mathbf{A})} \text{adj}(\mathbf{A})$ . As  $\mathbf{A}$  becomes more linearly dependent and  $\det(\mathbf{A}) \rightarrow 0$ ,  $\|\mathbf{A}^{-1}\|$  can become so large it causes numerical errors. Rewriting the original equation:

$$\begin{aligned}\mathbf{y} &= \mathbf{x}\mathbf{x}_D^{-1}\mathbf{y}_D \\ &= \frac{1}{\det(\mathbf{x}_D)} \mathbf{x} \text{adj}(\mathbf{x}_D) \mathbf{y}_D\end{aligned}$$

The errors caused by  $\det(\mathbf{x}_D) \rightarrow 0$  propagate to  $\mathbf{y}$ , causing predictions to be wildly inaccurate anywhere outside of the original training set.

#### Practical Implications

1. Instability: With a small determinant, minor variations in the attributes can lead to significant variations in the results. So, a small difference in income might result in a disproportionate change in loan eligibility.
2. Poor Generalization: If the matrix is based on data with limited variation (like our small community example), it's essentially trained on a very narrow subset of potential applicants. If someone from outside this narrow subset applies (e.g., a person with a 2-year employment but a \$70,000 income), the system may not process their application fairly or accurately because it's unfamiliar with such profiles.

**Given that a matrix used for determining loan approvals has a determinant close to zero due to limited variation in applicants' attributes:**

*Which of the following implications might this have on the decision-making process? Choose all options that apply.*

- A) It ensures a more uniform scoring system since most applicants have similar attributes.
- B) It can lead to unpredictable scores, where tiny variations in attributes yield vastly different outcomes.
- C) The system is more resilient to errors because of the limited attribute variation.
- D) It might not generalize well to broader populations, potentially leading to biases when applied to more diverse applicant groups.

Answer: B, D

B: When the determinant is close to 0, the denominator becomes too small, meaning a slight change can make the number unstable.

D: Since the matrix is trained with limited variation in applicants' attributes, it will not generalize to broader populations with vastly different attributes.

## 6.2 Voting and Probabilistic Models [5pts]

A country uses a probabilistic model to predict election outcomes and determine where resources (such as campaign funding or polling stations) should be allocated. The model uses factors like historical voting patterns, demographic data, voter turnout rates, and regional economic indicators to predict the probability of a particular candidate or party winning in each region.

However, the data used to train the model is incomplete: it primarily reflects urban areas, leaving rural voting behaviors underrepresented; some ethnic and socioeconomic groups have historically low voter turnout rates, meaning their data is sparsely included; and the model relies on historical data that may not accurately reflect recent political, economic, or social changes.

This results in several key issues. The model prioritizes campaign resources and polling stations in regions with high probabilities of voter influence, often favoring well-represented demographics, while rural or underrepresented regions may receive fewer polling stations, making it harder for individuals in those areas to vote. If polling stations are removed from low-priority areas due to the model's predictions, it could discourage voting in already marginalized communities, further disenfranchising groups with historically low turnout and perpetuating cycles of political exclusion.

**Which of the following is an ethical way to address the issue of bias in a probabilistic model used for allocating voting resources?**

- A) Collect and incorporate diverse and representative data to ensure the model accounts for voting patterns across all regions, demographics, and socioeconomic groups.
- B) Prioritize resource allocation only in regions with historically high voter turnout, as these areas are statistically more likely to impact election outcomes.
- C) Introduce fairness constraints in the model to guarantee equitable resource distribution, ensuring underrepresented regions and groups are not disadvantaged.
- D) Exclude regions with low voter turnout from the model's predictions, as their voting behavior is less predictable and may skew results.

Answer: A, C

A: Including all groups and regions will better represent the entire community with fairness.

C: Fairness constraints can ensure marginalized regions to get reasonable amount of resources.

## 7 Programming [5 pts]

See the Programming subfolder in Canvas.

## 8 Bonus Questions [7% Bonus for All]

### 8.1 Marginal Probability Density Functions (2% Bonus for All)

Suppose that  $X$  and  $Y$  have joint pdf given by

$$f_{X,Y}(x,y) = \begin{cases} 2e^{-2y}, & 0 \leq x \leq 1, y \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

What are the marginal probability density functions for  $X$  and  $Y$ ? [5 pts]

$$f_X(x) = \int_0^\infty f_{X,Y}(x,y) dy$$

$$= \int_0^\infty 2e^{-2y} dy$$

$$= [-e^{-2y}]_0^\infty = 1$$

$$f_Y(y) = \int_0^1 f_{X,Y}(x,y) dx$$

$$= \int_0^1 2e^{-2y} dx$$

$$= [2e^{-2y} \cdot x]_0^1 = 2e^{-2y}$$

$$\therefore f_X(x) = \begin{cases} 1, & 0 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

$$\therefore f_Y(y) = \begin{cases} 2e^{-2y}, & y \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

## 8.2 Coin Toss Game (2% Bonus for All)

A person decides to toss a biased coin with  $P(\text{heads}) = \frac{1}{3}$  repeatedly until he gets a head. He will make at most 6 tosses. Let the random variable  $Y$  denote the number of heads. Find the pmf of  $Y$ . Then, find the variance of  $Y$ . Round your answer to 3 decimal places. (*It is possible to thoroughly support your answer to this question in 5 to 10 lines*) [poin total here]

$$P(Y = k) = \begin{cases} \left(\frac{2}{3}\right)^6, & k = 0 \\ 1 - \left(\frac{2}{3}\right)^6, & k = 1 \end{cases}$$

$$P(Y = 0) \approx 0.088$$

$$P(Y = 1) \approx 0.912$$

$$E[Y] = 0 \cdot P(Y = 0) + 1 \cdot P(Y = 1) = 0.912$$

$$E[Y^2] = 0^2 \cdot P(Y = 0) + 1^2 \cdot P(Y = 1) = 0.912$$

$$\text{Var}(Y) = E[Y^2] - (E[Y])^2 = 0.912 - 0.912^2 = 0.080$$

### 8.3 Dice Roll Expectation (3% Bonus for All)

Suppose you roll an 8-sided die. For each roll:

- You are paid the face value of the roll.
- If the roll gives  $Y \in \{1, 2, 3, 4, 5\}$ , the game stops.
- If the roll gives  $Y \in \{6, 7, 8\}$ , you can roll again.

The probabilities are uniform, and the payoff structure is:

- For  $Y \in \{1, 2, 3, 4, 5\}$ , the expected value of the payoff is 3.
- For  $Y \in \{6, 7, 8\}$ , the expected value of the payoff is 7 plus the extra roll's expected value.

What is the expected payoff for this game? [5pts]

$$E[Y] = \frac{5}{8} \cdot 3 + \frac{3}{8} \cdot (7 + E[Y])$$

$$E[Y] = \frac{15}{8} + \frac{21}{8} + \frac{3}{8}E[Y]$$

$$\frac{5}{8}E[Y] = \frac{36}{8}$$

$$\therefore E[Y] = \frac{36}{5} = 7.2$$