

Data Science (CSE558)

End-Sem

Time: 2 Hours

Maximum Marks: 50

Instructions

1. Please read all questions carefully before writing your answers. The 1-mark questions need to answered with a single word. An incorrect answer to these questions fetches negative marks.
2. Please write neat, short and meaningful answers. Untidy answers will not be evaluated.
3. All parts of a single question should be answered collectively, i.e., all at one place.

All The Best!

1. Consider a weighted coin such that probability of getting head is 0.2.
 - a. What is the probability of getting at least 40 heads from 100 random tosses of this coin? (3 Marks)
 - b. Use this coin to define a random variable that can be used to design a family of random matrices of size $k \times d$ such that a matrix randomly sampled from the family ensure preservation of distance between 2^k pairs of points. (5 Marks)
2. Let A and B be two events in a sample space such that $\Pr(A) = 6/11$ and $\Pr(B) = 6/13$.
 - a. Can the events be disjoint? (+1/-0.5 Marks)
 - b. Prove or disprove your answer. (2 Marks)
3. Recall the hypothesis tests for population means, i.e., z-test and t-test. For a sample of size n, after computing the test statistics why degree of freedom is taken into account in t-test whereas not in z-test? (2 Marks)
4. Consider a bag with four unbiased dice as $\{d_1, d_2, d_3, d_4\}$. For every $1 \leq i \leq 4$, d_i has i-faces numbered 1 to i.
 - a. Let one die be selected uniformly at random from the bag and it has been rolled 3 times. The outcome is $\{1, 3, 1\}$. Use maximum likelihood estimation to determine the most likely dice selected from the bag. (2 Marks)
 - b. Use Bayes' rule to compute the probability of the most probable die for the given outcome after every roll. Compute $\Pr(\text{die}|\{1\})$ after 1st roll, compute $\Pr(\text{die}|\{3\})$ after 2nd roll and so on. (4 Marks)
 - c. Let a die be selected from the bag to play snakes and ladder between three friends. Each player can start her game only when she rolls 1. Let player-1, Player-2, and Player-3 take 3, 7, and 2 rolls to start their respective games. Using maximum likelihood estimation determine, the most likely die selected from the bag. (3 Marks)
5. A random point v on the surface of a unit sphere in a d -dimensional space, centered at origin is generated as follow: for every $i \in [1, \dots, d]$, $v(i)$ is either $+(d^{-0.5})$ or $-(d^{-0.5})$ with equal probability. Let x and y be two such points.
 - a. Prove that as d tends to infinity, the points tend to be orthogonal. (2 Marks)
 - b. For a fixed d , use Chebyshev's to bound the probability of the event $|x^T y| \geq 0.1$. (3 Marks)

- c. For a fixed d , use Bernstein to prove that there are at most $O(2^{0.01d})$ points such the following event E is true with probability at least 0.9.
E: For all pair of points (x, y) from the set, $|x^T y| < 0.1$. (5 Marks)
6. Let A be a matrix of size $n \times m$ such that its rank is 10, where $10 < \min(n, m)$. Let $A = U\Sigma V^T$ where U and V are two orthonormal square matrices of size $n \times n$ and $m \times m$ respectively.
 a. How many non-zero singular values are there in Σ ? (+1/-0.5 Marks)
 b. Let, $k < 10$ and $Z = [U_{11-k}, U_{12-k}, \dots, U_9, U_{10}]$ consisting of k singular vectors from U . $B = ZZ^T A$ be the low rank representation of A . Then calculate the exact Frobenius norm of the difference matrix between the original data and the projected data i.e., $\|A - B\|_F^2$. (3 Marks)
7. Continuing the above question, consider another matrix B of size $m \times n$. The running time of AB is $O(n^2m)$, which worries you because it has a quadratic dependence on n . So, your friend proposed the following steps to improve the running time of the matrix-matrix product.
 - Sample a random vector ' q ' in R^m such that, every index of q is -1 or +1 with equal probability.
 - Compute $C = Aq$ and $D = q^T B$.
 - Return $X = CD$.
 a. What is the running time of the above algorithm? (+1/-0.5 Marks)
 b. You noticed that, $E[qq^T] = I$ ($m \times m$ identity matrix). Hence, $E[X] = AB$ and it is an unbiased estimator. If the $\text{rank}(A) = \text{rank}(B) = \text{rank}(AB) = 10$ then what is the rank of X ? Why? (2 Marks)
8. Let A be a matrix of size $n \times m$. Let C be a set of $2k$ columns in A and R be set of k rows in A . With the selected C and R design an efficient algorithm to compute U for the CUR decomposition. State its running time? (5 Marks)
9. Let palyer1 and player2 are playing for car1 and car2 in two different Monty Hall game separately. Let player1 did not change her preference upon given a choice and she lost the car1. In a separate game let player2 decides to change her preference upon given the option. With this strategy is the player2 guaranteed to win the car2 in her game? (+1/-0.5 Marks)
10. Let H_m be a $n \times n$ Hadamard matrix, where $n = 2^m$ for some positive integer m . Let $H_m = \begin{pmatrix} H_{m-1} & H_{m-1} \\ H_{m-1} & -H_{m-1} \end{pmatrix}$ such that $H_0 = [1]$. Let v be a n -dimensional vector. Write an efficient pseudocode for the matrix vector product $H_m v$. What is its running time? (5 Marks)
11. (Bonus) Recall that sparsity in JL matrix has an advantage in the running time. Is this always useful? Explain your reason with an analysis. (3 Marks)
-

Rubric

Q1.

Q1 (a) Random Variable x

$$x = \begin{cases} 1 & \text{we assign} \\ & 0.2 \\ 0 & \text{else} \end{cases}$$

$$100 \cdot E[x] = 100 \cdot (1 \times 0.2 + 0 \times 0.3) \\ = 20 \text{ expected heads.}$$

0.5
Mark

$$R = \sum_{i=1}^{100} x_i$$

$$E[R] = 20$$

Markov

we look at 40

$$Pr(R \geq 40) \leq \frac{E[R]}{40} = \frac{1}{2}$$

Markov

1
Mark

$$\text{Variance}(R) = \sum_{i=1}^{100} \text{Var}(x_i) = 100 \cdot \text{Var}(x) \\ = 100[E[x^2] - [E[x]]^2]$$

$$= 100[0.2 - 0.04] = 16$$

0.5
Mark

0.5
Chebyshev

$$Pr(|R - E[R]| \geq 20) \leq \frac{\text{Var}(R)}{(20)^2} = \frac{16}{20 \cdot 20} = \frac{1}{25}$$

0.5
Mark

of heads ≥ 40

$$\text{Prob} \leq \frac{1}{25} \quad \text{Full Marks}$$

OR

$$\sum_{i=40}^{100} \binom{100}{i} \cdot i^{0.2} \cdot 60^{0.8} \quad \left. \vphantom{\sum_{i=40}^{100}} \right\} 0.5 \text{ Mark}$$

(b) $E[x] = 0 \rightarrow 1 \text{ Mark}$
 $Var(x) = 1/k \rightarrow 2 \text{ Marks}$

Unbiased coin

HT $\rightarrow \tilde{H}$	\rightarrow Prob	0.16
TH $\rightarrow \tilde{T}$	\rightarrow Prob	0.16

$\leftarrow \Omega$ new sample space

$\times \begin{cases} TT & \rightarrow \text{" } 0.64 \\ HH & \rightarrow \text{" } 0.04 \end{cases}$

$$P_{\tilde{H}}(\tilde{H}) = \frac{0.16}{0.32} = 0.5$$

$$P_{\tilde{T}}(\tilde{T}) = 0.5$$

} 1 Mark

$$x = \begin{cases} \frac{1}{\sqrt{k}} & \text{if } \tilde{H} \\ -\frac{1}{\sqrt{k}} & \text{if } \tilde{T} \end{cases}$$

} 1 Mark

Prove:-

$$E[x] = 0$$

$$Var(x) = \frac{1}{k}(0.16) + \frac{1}{k}(0.16) - (E[x])^2$$

we need 0

Prove:-

$$\frac{0.5}{k} + \frac{0.5}{k}$$

$$Var = \frac{1}{k}$$

Que.3

1. For Z-test - No degree of freedom because population variance is known and fixed.
2. For T-test - DoF is critical because the test relies on sample variance and smaller sample size increases uncertainty. **2 Marks if the above points are mentioned in the answer.**

Que.7

Time complexity = $O(n^2 + nm)$ ----- **1 M / -0.5M**

Rank of X = 1 ----- **1 M**

Reason(C, D only 1,1 vectors) ----- **1 M**

Que.8

SVD method :

Running time — **1M** Running time = $O(K^3)$

Intersection — **2M**

Pseudo Inverse — **2M**

Other method :

Calculation of U — **1M**

C+, R+ calculation — **2M**

Efficient Algorithm ----- **1M**

Time Complexity ----- **1M** Running time = $O(k^2n + kmn + k^2m)$

Marks awarded for both methods

Que.9

Answer : **No, 1 Mark for correct, -0.5 for incorrect.**

Que-2

Q2.

a) No

$$b) P(A) = \frac{6}{11} = \frac{6 \times 13}{11 \times 13} = \frac{78}{143}$$

$$P(B) = \frac{6}{13} = \frac{6 \times 11}{13 \times 11} = \frac{66}{143}$$

If disjoint then A & B cannot have a common sample --- (i)

If (i) true then $P(A) + P(B) \leq 1$ --- (ii)

$$P(A) + P(B) = \frac{144}{143} > 1$$

\therefore Contradiction

Que.4

Q4/ Given :-

$$d_1 = \{1\}$$

$$d_2 = \{1, 2\}$$

$$d_3 = \{1, 2, 3\}$$

$$d_4 = \{1, 2, 3, 4\}$$

- Each die has is equally likely to be chosen from the bag initially, so the prior probability of selecting

a) Maximum likelihood Estimation for outcome $\{1, 3, 1\}$

$$d_1: P(\{1, 3, 1\} | d_1) = 1 \times 0 \times 1 = 0$$

$$d_2: P(\{1, 3, 1\} | d_2) = \frac{1}{2} \times 0 \times \frac{1}{2} = 0$$

$$d_3: P(\{1, 3, 1\} | d_3) = \left(\frac{1}{3}\right) \left(\frac{1}{3}\right) \left(\frac{1}{3}\right) = \frac{1}{27}$$

$$d_4: P(\{1, 3, 1\} | d_4) = \left(\frac{1}{4}\right) \left(\frac{1}{4}\right) \left(\frac{1}{4}\right) = \frac{1}{64}$$

Ans: d_3

• Ans = d_3 (involves conditional probability) 2 Marks
 (or) if conditional probability - 3 Marks

(i) Geometric distribution formula is used and correct answer - 3 Marks

(ii) eqn ① is written and substitution of $\theta = 1, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}$ is done - 3 Marks

(iii) directly d_4 is written as answer - 1 Mark

(iv) if $\theta = \frac{n}{\sum x_i} = \frac{3}{12} = \frac{1}{4}$ is computed only 2.5 Marks

$$= \frac{1 \cdot \frac{1}{4}}{(1 \cdot \frac{1}{4}) + (\frac{1}{2} \cdot \frac{1}{4}) + (\frac{1}{3} \cdot \frac{1}{4}) + (\frac{1}{4} \cdot \frac{1}{4})}$$

$$P(d_1|\{1\}) = \frac{1}{1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4}} = 0.48, \quad P(d_2|\{1\}) = \frac{1}{2}$$

$$P(d_3|\{1\}) = \frac{1}{3}, \quad P(d_4|\{1\}) = \frac{1}{4}$$

$$= 0.16, \quad = 0.12$$

After the first roll (1), the posterior probabilities are:-

- $P(d_1|\{1\}) = 0.48$
- $P(d_2|\{1\}) = 0.24$
- $P(d_3|\{1\}) = 0.16$
- $P(d_4|\{1\}) = 0.12$

After the second roll (Outcome = 3)

$$P(d_1|\{1,3\}) = \frac{P(3|d_1) \cdot P(d_1)}{P(3)} = \frac{0 \cdot (\frac{1}{4 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4}})}{(\quad)}$$

$$P(d_1|3) = 0$$

$$P(d_2|3) = 0 \quad (\text{since } P(3|d_2) = 0)$$

$$P(d_3|3) = \frac{P(3|d_3) \cdot P(d_3)}{P(3)}$$

$$= \frac{\frac{1}{3} \cdot \frac{1}{3}}{\frac{1}{3} \cdot (\frac{1}{3}) + \frac{1}{4} \cdot \frac{1}{4}}$$

$$P(d_3|3) = \frac{16}{25}$$

$$P(d_4|3) = \frac{P(3|d_4) \cdot P(d_4)}{P(3)}$$

$$= \frac{\frac{1}{4} \cdot \frac{1}{4}}{\frac{1}{3} \cdot (\frac{1}{3}) + \frac{1}{4} \cdot \frac{1}{4}}$$

After the second roll (3), the posterior probabilities are:-

- $P(d_1|\{1,3\}) = 0$
- $P(d_2|\{1,3\}) = 0$
- $P(d_3|\{1,3\}) = \frac{16}{25}$
- $P(d_4|\{1,3\}) = \frac{9}{25}$

$$P(d_3|1) =$$

After the third roll (outcome = 1):

$$P(d_3|1) = \frac{P(1|d_3) \cdot P(d_3)}{\frac{1}{3} \cdot \frac{16}{25} + \frac{1}{4} \cdot \frac{9}{25}} = \frac{64}{91}$$

$$P(d_4|1) = \frac{\frac{1}{4} \cdot \frac{9}{25}}{\frac{1}{3} \cdot \frac{16}{25} + \frac{1}{4} \cdot \frac{9}{25}} = \frac{27}{91}$$

After the third roll (1), the posterior probabilities are:

$$\bullet P(d_3|\{1, 3, 13\}) = 64/91$$

$$\bullet P(d_4|\{1, 3, 13\}) = 27/91$$

(i) all validations + correct answer - 4 marks
are shown

(or)

(ii) only d_3 calculation is done completely - 1 mark

(or)

(iii) Iteration 1 only - 1 mark

" (1+2) - 2 marks

" (1+2+3) - 4 marks
+ (updated probability)

Iteration (1 and 2) - 3 marks
+ (updated probability)

Que.5

Q5. $i \in \{1, 2, \dots, d\}$
 $v_i \in \left\{ \frac{1}{\sqrt{d}}, -\frac{1}{\sqrt{d}} \right\}$
 $x, y \in \mathbb{R}^d$

(a)

$$E[x^T y] = E\left[\sum_{i=1}^d x_i y_i\right] = \sum_{i=1}^d E[x_i y_i] \quad \left. \begin{array}{l} \text{Since } x \text{ \& } y \text{ are independent.} \\ \end{array} \right\} 0.5$$

$$E[x_i y_i] = E[x_i] \cdot E[y_i]$$

$$x_i, y_i = \begin{cases} 1/\sqrt{d} & p=1/2 \\ -1/\sqrt{d} & 1/2 \end{cases} \Rightarrow E[x_i] = E[y_i] = 0 \quad \left. \right\} 0.5$$

$$\Rightarrow E[x^T y] = \sum_{i=1}^d E[x_i y_i] = \sum_{i=1}^d E[x_i] E[y_i] = 0 \quad \left. \right\} 0.5$$

$$\begin{aligned} \text{Var}[x^T y] &= \text{Var}\left(\sum_{i=1}^d x_i y_i\right) = \sum_{i=1}^d \text{Var}(x_i y_i) \\ &= \sum_{i=1}^d \left(E[(x_i y_i)^2] - (E[x_i y_i])^2 \right) \\ &= \sum_{i=1}^d \left(E[x_i^2] \cdot E[y_i^2] - 0 \right) \\ &= \sum_{i=1}^d \frac{1}{d} \times \frac{1}{d} = \frac{1}{d} \quad \left(E[x_i^2] = \frac{1}{d} \right) \end{aligned} \quad \left. \right\} 1$$

$$\text{as } d \rightarrow \infty \quad \text{Var}[x^T y] = \lim_{d \rightarrow \infty} \frac{1}{d} \rightarrow 0$$

$\text{Var}[x^T y] \rightarrow 0$ implies $x^T y$ approaches one point value
 and as $E[x^T y] = 0$, this point is 0 .
 $\Rightarrow x^T y \rightarrow 0$ for $d \rightarrow \infty$ i.e. x, y are orthogonal

(b) Let Z be the R.V. for representing $x^T y$

$$\mu = E[x^T y] = 0$$

$$\text{Var}(Z) = \text{Var}(x^T y) = 1/d = \sigma^2$$

Chebyshev's Inequality:

$$P(|Z - \mu| > k\sigma) \leq \frac{1}{k^2}$$

$$\Rightarrow P(|Z - 0| \geq 0.1) \leq \frac{\sigma^2}{(0.1)^2} \leq \frac{1}{d^2 \cdot 0.1^2}$$

$$= \frac{100}{d^2} \cdot 0.5 \quad \left(k = \frac{0.1}{\sigma} \right)$$

(c) Bernstein's Inequality:

$$P(|S - E[S]| \geq t) \leq \exp\left(-\frac{t^2}{\sigma^2 + bt}\right) \quad \left. \vphantom{\exp}\right\} 0.5$$

$$X_i = x_i^T y_i = \begin{cases} 1/d & p = 1/2 \\ -1/d & b = 1/2 \end{cases} \quad \left. \vphantom{\begin{cases} 1/d \\ -1/d \end{cases}}\right\} 0.5$$

$$S = \sum_{i=1}^d X_i, \quad (X_i - E[X_i] = |x_i^T y_i - 0|) \quad \left. \vphantom{\sum}\right\} 0.5$$

$$E[S] = E[x^T y] = 0 \quad // \quad = |x_i^T y_i| = \frac{1}{d}$$

$$\Rightarrow b = 1/d$$

$$\sigma^2 = \text{Var}(X_i) = \text{Var}(x^T y) = 1/d \quad \left. \vphantom{\text{Var}}\right\} 0.5$$

$$|x^T y| \geq t = 0.1 \quad \left. \vphantom{|x^T y|}\right\} 0.5$$

Using the inequality

$$\Rightarrow P(|S - 0| \geq 0.1) \leq \exp\left(-\frac{(0.1)^2}{\frac{1}{d} + \frac{0.1}{d}}\right) \quad \left. \vphantom{\exp}\right\} 0.5$$

$$\Rightarrow P\left(\left|\sum_{i=1}^d x_i^T y_i\right| \geq 0.1\right) \leq \exp\left(-\frac{d}{110}\right)$$

$$\Rightarrow P(|x^T y| \geq 0.1) = 1 - P(|x^T y| < 0.1) \quad \left. \vphantom{P}\right\} 0.5$$

$$\leq \exp\left(-\frac{d}{110}\right)$$

$$\Rightarrow P(|x^T y| < 0.1) \geq 1 - \exp\left(-\frac{d}{110}\right)$$

If there are m points which form the set,

$$\text{no. of pairs } (x, y) = \frac{m(m-1)}{2}$$

$$\begin{aligned} P \left(P(|x^T y| < 0.1) \text{ for all pairs} \right) \\ = \left(1 - e^{-d/10} \right)^{\frac{m(m-1)}{2}} > 0.9 \end{aligned} \quad \left. \vphantom{\begin{aligned} P \left(P(|x^T y| < 0.1) \text{ for all pairs} \right) \\ = \left(1 - e^{-d/10} \right)^{\frac{m(m-1)}{2}} > 0.9 \end{aligned}} \right\} 1$$

$$\text{as } e^{-d/10} \rightarrow 0 < 1$$

$$\Rightarrow 1 - \frac{m(m-1)}{2} e^{-d/10} > 0.9$$

$$\Rightarrow \frac{m(m-1)}{2} e^{-d/10} \leq \frac{(0.1)2}{m(m-1)}$$

$$\Rightarrow m(m-1) \leq (0.2) e^{d/10}$$

$$\Rightarrow m \leq O \left(e^{d/20} \right) \leq O \left(2^{d/100} \right) \quad \left. \vphantom{\Rightarrow m \leq O \left(e^{d/20} \right) \leq O \left(2^{d/100} \right)} \right\} 0.5$$

Que.6

6.a: **10, 1 Mark** for correct, **-0.5** for incorrect.

6.b: **1 Mark** for expression B, **2 marks** for exact frobenius norm.

ZZ^T is projection Matrix which is spanned by Z , $Z = [U_{1-k}, U_{2-k}, \dots, U_{10}]$.

and, All the vectors in 'U' matrix is Orthonormal i.e. every pair wise vector is zero. that, effectively

$$B = \sigma_8 U_8 V_8^T + \sigma_9 U_9 V_9^T + \sigma_{10} U_{10} V_{10}^T \text{ --- } \underline{1 \text{ Mark}}$$

Assuming $k=3$.

After $\|A-B\|_F^2 =$

$$\left\| \sum_{i=1}^{10} \sigma_i U_i V_i^T - \sum_{j=8}^{10} \sigma_j V_j^T U_j \right\|_F^2$$

$\therefore U_i V_j$ are orthonormal from 10 to 8 get Cancel out to zero

$$\therefore \left\| \sum_{i=1}^7 \sigma_i U_i V_i^T \right\|_F^2$$

$$\Rightarrow \sum_{i=1}^7 \sigma_i^2 \text{ --- } \underline{2 \text{ Mark}}$$

Que-10

Q10

Base case: $H_0 = [1]$

$$H_1 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

and so on

for a vector $\vec{v}_m \rightarrow \begin{bmatrix} \vec{v}_1 \\ \vec{v}_2 \end{bmatrix} \quad \forall n \geq 1$

this split is always true as $n = 2^m$ & $m \geq 0$

until $m=0$ where $\vec{v} = [\vec{v}_0]$

for a general $\vec{v}_m = \begin{bmatrix} \vec{v}_1 \\ \vec{v}_2 \end{bmatrix}$ where $n = 2^m$

& a general

$$H_m = \begin{bmatrix} H_{m-1} & H_{m-1} \\ H_{m-1} & -H_{m-1} \end{bmatrix}$$

$$H_m \vec{v} = \begin{bmatrix} H_{m-1} & H_{m-1} \\ H_{m-1} & -H_{m-1} \end{bmatrix} \begin{bmatrix} \vec{v}_1 \\ \vec{v}_2 \end{bmatrix} = \begin{bmatrix} H_{m-1} \vec{v}_1 + H_{m-1} \vec{v}_2 \\ H_{m-1} \vec{v}_1 - H_{m-1} \vec{v}_2 \end{bmatrix}$$

~~divide each of the~~

→ Divide each of the 2 terms further until H_0 & \vec{v} length = 1 ($\log n$ operations)

→ Perform addition backwards from $m=0$ to $m=2^n$ (n operations on merge step/compute step)

* Idea similar to merge sort & can cache at intermediate steps for future calculations

Time Complexity: $n \log(n)$

Note: In Q10. of the endsem. If you have used a recursive algorithm. Partial marks have been awarded for correct divide and merge steps. And another mark for the correct time complexity mentioned. There are no marks for a rank approximation or any other randomized prediction of the result, an exact computation was required.

