

CS 771A: Intro to Machine Learning, IIT Kanpur			Endsem Exam (22 Nov 2022)	
Name	JAYA GUPTA			50 marks Page 1 of 6
Roll No	200471	Dept.	CSE	

### Instructions:

1. This question paper contains 3 pages (6 sides of paper). Please verify.
2. Write your name, roll number, department in **block letters** with **ink** on **each page**.
3. Write your final answers neatly **with a blue/black pen**. Pencil marks may get smudged.
4. Don't overwrite/scratch answers especially in MCQ.



**Q1. Write T or F for True/False in the box and give justification below. (4 x (1+2) = 12 marks)**

1	The Nikola company shares have a 40% chance of crashing if its owner Ksümnöle tweets something silly. The shares have a 10% chance of crashing if no silly tweet is sent. Ksümnöle tweets something silly with a 20% chance. Then, the probability that Nikola shares will crash, is less than 20%. Justify by calculating the probability.	T
$P[\text{Nikola shares crash}] = P[\text{Ksümnöle tweets silly}] \times P[\text{crash}   \text{silly tweet}] + P[\text{no silly tweet}] \times P[\text{crash}   \text{no silly tweet}]$ $= 0.2 \times 0.4 + 0.8 \times 0.1$ $= 0.08 + 0.08 = 0.16 = 16\% < 20\%$		
2	Given three vectors $x, y, z \in \mathbb{R}^2$ such that $x^T z > x^T y$ , it is always the case that $\ x - z\ _2^2 < \ x - y\ _2^2$ . Give a proof if True else give a counter example.	F
<p>Let <math>x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}</math> <math>y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}</math> <math>z = \begin{bmatrix} 1 \\ 7 \end{bmatrix}</math> <math>x^T z = 1 &gt; x^T y = 0</math></p> <p><math>\ x - z\ _2^2 = \left\  \begin{bmatrix} 0 \\ -7 \end{bmatrix} \right\ _2^2 = 49</math> <math>\ x - y\ _2^2 = \left\  \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right\ _2^2 = 2</math></p> <p>Clearly <math>\ x - z\ _2^2 &gt; \ x - y\ _2^2</math>. Hence above st. is false.</p>		
3	Consider the set $\mathcal{X} = \{-1, +1\}^3$ of 3D vectors with $\pm 1$ coordinates. Any map $\phi: \mathcal{X} \rightarrow \mathbb{R}^d$ s.t. for all $x, y \in \mathcal{X}$ , $\phi(x)^T \phi(y) = (1 + x^T y)^2$ must use $d \geq 10$ dims. Give a proof if True else give a map using fewer dimensions as a counter example.	F (false)
<p><math>(1 + x^T y)</math> is quadratic kernel <math>= (1 + \langle x, y \rangle)^2 = 1 + 2\langle x, y \rangle + \langle x, y \rangle^2</math></p> <p><math>= 1 + 2\langle x, y \rangle + \sum_{j=1}^3 \sum_{i=1}^3 x_i x_j y_i y_j</math></p> <p>So <math>\phi(x) = [\phi_1(x) \ \phi_2(x) \ \phi_3(x)]</math></p> <p><math>\phi_1(x) = 1 \rightarrow (1)</math></p> <p><math>\phi_2(x) = [\sqrt{2} x] \rightarrow d = 3</math></p> <p><math>\phi_3(x) = [x_i x_j]_{i,j \in \{1,2,3\}} \rightarrow d^2 = 9 \rightarrow</math> Reduce to <math>d = 6</math>. as <math>\{x_1 x_1, x_1 x_2, x_1 x_3, x_2 x_1, x_2 x_2, x_2 x_3, x_3 x_1, x_3 x_2, x_3 x_3\}</math></p> <p>Further <math>\{x_1 x_1, x_2 x_2, x_3 x_3, x_1 x_2, x_2 x_1, x_1 x_3, x_3 x_1, x_2 x_3, x_3 x_2\}</math></p> <p><math>\Rightarrow d + d + 1 = 3 + 3 + 1 = 7</math></p> <p><math>\phi(x)</math> should have at least 7 dims.</p>		

Since  $x_i = \pm 1, x_i^2 = 1$  Further  $\{x_1 x_1, x_2 x_2, x_3 x_3, x_1 x_2, x_2 x_1, x_1 x_3, x_3 x_1, x_2 x_3, x_3 x_2\}$

$\Rightarrow$  Reduce to  $d = 4$



- 4 If  $X, Y \in \mathbb{R}^{3 \times 3}$  are rank one matrices, then  $X + Y$  can never be rank one, no matter what are  $X, Y$ . Give a brief proof if True else give a counter example.

F

$$X = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix} \quad Y = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad X+Y = \begin{bmatrix} 2 & 4 & 6 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix}$$

rank=1                      rank=1                      rank=1.

$C_2 = 2C_1 \rightarrow (\text{column 1})$   
 $C_3 = 3C_1$

**Q2. (Informative non-response models)** Melbo is studying how one's income level affects one's reluctance to reveal one's income publicly.  $n$  people were chosen with incomes  $X_1, X_2, \dots, X_n$ . Melbo knows that the income levels  $X_i$  are distributed as independent standard Gaussian random variables i.e.,  $X_i \sim \mathcal{N}(0,1)$  for all  $i$  (let us interpret positive  $X_i$  as higher-than-median income and negative  $X_i$  as lower-than-median income). However, not everyone wants to reveal their income. When Melbo conducts the survey, the responses are  $Z_1, Z_2, \dots, Z_n$ . If the  $i^{\text{th}}$  person reveals their income, then  $Z_i = X_i$  else  $Z_i = \phi$ . It is known that  $\mathbb{P}[Z_i \neq \phi \mid X_i] = \exp\left(-\frac{\alpha^2 X_i^2}{2}\right)$ , where  $\alpha > 0$  is an unknown parameter to be learnt. **(Total 12 marks)**

1. Is a rich person e.g.,  $X_i = 100$  more likely or less likely to reveal their income than a person with close-to-median income e.g.,  $X_j = -0.01$ ? Give brief justification. (1+1 = 2 marks)

We need to compare  $\mathbb{P}[Z_i \neq \phi \mid X_i]$ .

Rich person:-  $\mathbb{P}[Z_i \neq \phi \mid X_i = 100] = \exp\left(-\frac{\alpha^2 10000}{2}\right) = \exp(-5000\alpha^2)$

Median person:-  $\mathbb{P}[Z_j \neq \phi \mid X_j = -0.01] = \exp\left(-\frac{\alpha^2 0.0001}{2}\right) = \exp(-0.00005\alpha^2)$

$e^{-x^2}$  is decreasing func, hence close to median person is more likely to reveal their income than rich person.  $\rightarrow$  less likely.

2. Is a poor person e.g.,  $X_i = -10$  more likely or less likely to reveal their income than a person with close-to-median income e.g.,  $X_j = 0.1$ ? Give brief justification. (1+1 = 2 marks)

Poor person:-  $\mathbb{P}[Z_i \neq \phi \mid X_i = -10] = \exp\left(-\frac{100\alpha^2}{2}\right) = \exp(-50\alpha^2)$

Median person:-  $\mathbb{P}[Z_j \neq \phi \mid X_j = 0.1] = \exp\left(-\frac{0.01\alpha^2}{2}\right) = \exp(-0.005\alpha^2)$

So,  $\mathbb{P}[Z_i \neq \phi \mid X_i = -10] < \mathbb{P}[Z_j \neq \phi \mid X_j = 0.1]$

Hence poor person is less likely to reveal their income than close-to-median income person.

3. Derive an expression for  $\mathbb{P}[Z_i \neq \phi]$  the prior probability of a person revealing their income. Show steps and give your answer as a function  $h(\alpha)$ . **Hint:** the density of a Gaussian looks like  $\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$  and  $X_i \sim \mathcal{N}(0,1)$ . Also,  $\int_{-\infty}^{\infty} \exp\left(-\frac{a^2 t^2}{2}\right) dt = \sqrt{\frac{2\pi}{a^2}}$ . (4 marks)

$$\text{like } \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \text{ and } X_i \sim \mathcal{N}(0,1). \text{ Also, } \int_{-\infty}^{\infty} \exp\left(-\frac{a^2 t^2}{2}\right) dt = \sqrt{\frac{2\pi}{a^2}}. \quad (4 \text{ marks})$$

CS 771A: Intro to Machine Learning, IIT Kanpur			Endsem Exam (22 Nov 2022)	
Name	JAYA GUPTA			50 marks Page 3 of 6
Roll No	200421	Dept.	CSE.	

$$\begin{aligned}
 P[Z_i \neq \phi] &= \int_{x_i \in \mathcal{R}} P[Z_i \neq \phi, x_i] dx_i = \int_{-\infty}^{\infty} P[Z_i \neq \phi | x_i] P[x_i] dx_i \\
 &\quad \text{law of total prob.} \quad P(A|B) = \frac{P(A, B)}{P(B)} \\
 &= \int_{-\infty}^{\infty} \exp\left(-\frac{\alpha^2 x_i^2}{2}\right) \underbrace{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x_i^2}{2}\right)}_{\mathcal{N}(0,1) \sim x_i} dx_i \\
 &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \exp\left(-(\alpha^2 + 1) \frac{x_i^2}{2}\right) dx_i \quad \left[ \text{given: } \int_{-\infty}^{\infty} \exp\left(-\frac{\alpha^2 t^2}{2}\right) dt = \sqrt{\frac{2\pi}{\alpha^2}} \right] \\
 &= \frac{1}{\sqrt{2\pi}} \times \sqrt{\frac{2\pi}{\alpha^2 + 1}} = \frac{1}{\sqrt{\alpha^2 + 1}}
 \end{aligned}$$

$$h(\alpha) = \frac{1}{\sqrt{\alpha^2 + 1}} = P[Z_i \neq \phi]$$

4. Write down an expression for the negative log-likelihood of the form (no derivation needed)

$$\begin{aligned}
 &P[Z_i \neq \phi, x_i] = P[Z_i \neq \phi | x_i] P[x_i] \\
 \mathcal{L}(\alpha) &= - \sum_{i: Z_i \neq \phi} \ln P[Z_i \neq \phi, x_i] - \sum_{i: Z_i = \phi} \ln P[Z_i = \phi]
 \end{aligned}$$

Notice that the terms in the first summation involve joint probability. (2 marks)

$$\mathcal{L}(\alpha) = \sum_{i: Z_i \neq \phi} \left( -\ln \frac{1}{\sqrt{2\pi}} + \frac{\alpha^2 x_i^2}{2} + \frac{x_i^2}{2} \right) - \sum_{i: Z_i = \phi} \ln \left( 1 - \frac{1}{\sqrt{\alpha^2 + 1}} \right)$$

$\underbrace{\frac{1}{\sqrt{\alpha^2 + 1}}}_{h(\alpha)}$

5. Write down an expression for the gradient  $\mathcal{L}'(\alpha)$  (no derivation needed). (2 marks)

$$\mathcal{L}'(\alpha) = \sum_{i: Z_i \neq \phi} \alpha x_i^2 - \sum_{i: Z_i = \phi} \frac{(\sqrt{\alpha^2 + 1} + 1)}{\alpha(\alpha^2 + 1)}$$



**Q3. (Quantile regression)** Can we find the  $k^{\text{th}}$  largest number in a set of  $n$  numbers simply by solving an optimization problem?! Turns out it is indeed possible using a trick called quantile regression. For a set of real numbers  $x_1 < x_2 < \dots < x_n$  (sorted in ascending order for sake of simplicity), for any integer  $k = 0, 1, 2, \dots, n$ , consider the problem  $\operatorname{argmin}_{z \in [x_1, x_n]} f_k(z)$ , with

$$f_k(z) \stackrel{\text{def}}{=} \left(\frac{k}{n} - 1\right) \cdot \sum_{x_i < z} (x_i - z) + \frac{k}{n} \cdot \sum_{x_i \geq z} (x_i - z)$$

There are no duplicates in  $x_1, \dots, x_n$ . Assume that an empty sum equals 0.

1. Find a minimizer for  $\operatorname{argmin}_{z \in [x_1, x_n]} f_n(z)$  i.e.,  $k = n$ . Show brief derivation. (1+1=2 marks)

$f_n(z) = \sum_{x_i \geq z} (x_i - z)$

Notice that  $f_n(z) \geq 0$  because  $(x_i - z) \geq 0$  as  $x_i \geq z$ .  
 Hence min value of  $f_n(z) = 0$ .

$\operatorname{argmin}_{z \in [x_1, x_n]} \sum_{x_i \geq z} (x_i - z) = \sum_{x_i \geq x_n} (x_i - z) = (x_n - z) = 0$

Min value of  $f_n(z)$  will be 0, when  $z = x_n$ .  
 Hence minimizer is  $\boxed{z = x_n}$ .

2. Find a minimizer for  $\operatorname{argmin}_{z \in [x_1, x_n]} f_0(z)$  i.e.,  $k = 0$ . Show brief derivation. (1+1=2 marks)

$f_0(z) = \sum_{x_i < z} (z - x_i)$

Notice that  $f_0(z) \geq 0$  because  $(z - x_i) \geq 0$  as  $x_i < z$ .  
 Hence min value of  $f_0(z) = 0$ .

$\operatorname{argmin}_{z \in [x_1, x_n]} \sum_{x_i < z} (z - x_i) = \sum_{x_i < x_1} (z - x_i) = \text{empty sum} = 0$

Min value of  $f_0(z)$  will be 0, when  $z = x_1$ .  
 Hence minimizer is  $\boxed{z = x_1}$ .

3. Let us handle  $k \in [1, n-1]$ . Show brief derivation that if  $x_j < a < b \leq x_{j+1}$ ,  $a \neq b$ , then

- We have  $f_k(a) > f_k(b)$  if  $1 \leq j < k$ .
- We have  $f_k(a) < f_k(b)$  if  $k < j < n$ , we have.
- We have  $f_k(a) = f_k(b)$  if  $j = k$ , i.e., for  $x_k < a < b \leq x_{k+1}$ . (4+4+4 = 12 marks)

After establishing a few more results like the ones above (which you do not have to show), we can deduce that any value of  $z \in [x_k, x_{k+1})$  is a minimizer of  $\operatorname{argmin}_{z \in [x_1, x_n]} f_k(z)$ . (Total 16 marks)

Name

JAYA GUPTA

50 marks

Roll No

200471

Dept.

CSE

Page 5 of 6

lets find  $f_k(a) - f_k(b)$ ,

$$f_k(a) = \left(\frac{k}{n} - 1\right) \sum_{x_i < a} (x_i - a) + \frac{k}{n} \sum_{x_i \geq a} (x_i - a)$$

$$f_k(b) = \left(\frac{k}{n} - 1\right) \sum_{x_i < b} (x_i - b) + \frac{k}{n} \sum_{x_i \geq b} (x_i - b)$$

Now since  $x_j < a < b \leq x_{j+1}$ ,  $a$  and  $b$  both lie b/w  $x_j$  and  $x_{j+1}$ , num of elements ( $x_i$ 's) s.t.  $x_i < a$  or  $x_i < b = j = \{x_1, x_2, \dots, x_j\}$ .

similarly num of elements s.t.  $x_i \geq a$  or  $x_i \geq b = (n-j) = \{x_{j+1}, \dots, x_n\}$ .

$$\text{so } f_k(a) - f_k(b) = \left(\frac{k}{n} - 1\right) \sum_{\substack{x_i < a \\ \text{or } x_i < b}} (b-a) + \frac{k}{n} \sum_{\substack{x_i \geq a \\ \text{or } x_i \geq b}} (b-a)$$

$$\Rightarrow (b-a) \left[ \left(\frac{k}{n} - 1\right) j + \frac{k}{n} (n-j) \right] = (b-a) (k-j)$$

$$a) f_k(a) > f_k(b) \Rightarrow f_k(a) - f_k(b) > 0$$

$$(b-a) (k-j) > 0$$

$$k > j \quad \text{Hence } 1 \leq j < k$$

Note:-  
 $1 \leq j < n$   
 $j \in \{1, 2, \dots, n-1\}$

$$b) f_k(a) < f_k(b) \Rightarrow f_k(a) - f_k(b) < 0 \Rightarrow (b-a) (k-j) < 0$$

$$\text{Hence } k < j < n$$

$$c) f_k(a) = f_k(b) \Rightarrow f_k(a) - f_k(b) = 0 \Rightarrow (b-a) (k-j) = 0$$

$$\text{Hence } k = j$$

**Q4. (Robust mean estimation)** Melbo has got samples  $X_1, \dots, X_n$  from a Gaussian with unknown mean  $\mu$  but known variance  $\sigma = \frac{1}{\sqrt{2\pi}}$  i.e., with density  $f(X; \mu) = \exp(-\pi(X - \mu)^2)$ . Melbo wishes to estimate  $\mu$  using these samples but is stuck since some samples were corrupted by Melbo's enemy Oblem. It is not known which samples did Oblem corrupt. Let's use latent variables to solve



this problem. For each  $i$ , we say  $Z_i = 1$  if we think  $X_i$  is corrupted else  $Z_i = 0$ . For any  $\mu \in \mathbb{R}$ , we are told that  $\mathbb{P}[Z_i = 1 | \mu] = \eta$ , and that  $\mathbb{P}[X_i | \mu, Z_i = 1] = \epsilon$ , and  $\mathbb{P}[X_i | \mu, Z_i = 0] = f(X_i; \mu)$ . Thus, we suspect that Oblem corrupted around  $\eta$  fraction of the samples and we assume that a corrupted sample can take any value with probability  $\epsilon$ . Assume  $\epsilon, \eta < \frac{1}{10}$  and are both known.

1. For a given  $\mu$ , derive for a rule to find out if  $\mathbb{P}[Z_i = 1 | X_i, \mu] > \mathbb{P}[Z_i = 0 | X_i, \mu]$  or not.

$$\mathbb{P}[Z_i = 1 | X_i, \mu] = \frac{\mathbb{P}[X_i | Z_i = 1, \mu] \cdot \mathbb{P}[Z_i = 1 | \mu]}{\mathbb{P}[X_i | \mu]} \quad (\text{Bayes' theorem})$$

$$\mathbb{P}[Z_i = 1 | X_i, \mu] > \mathbb{P}[Z_i = 0 | X_i, \mu]$$

$$\Rightarrow \mathbb{P}[X_i | Z_i = 1, \mu] \mathbb{P}[Z_i = 1 | \mu] > \mathbb{P}[X_i | Z_i = 0, \mu] \mathbb{P}[Z_i = 0 | \mu]$$

$$\Rightarrow \epsilon \eta > f(X_i; \mu) (1 - \eta) \quad \left[ \mathbb{P}[Z_i = 0 | \mu] = 1 - \eta \right]$$

$$\Rightarrow \epsilon \eta > \exp(-\pi (X_i - \mu)^2) (1 - \eta) \quad \left[ \text{If sign} = 1, \right.$$

$$\Rightarrow \frac{\epsilon \eta \exp(\pi (X_i - \mu)^2)}{1 - \eta} - 1 > 0. \quad \left. \begin{array}{l} Z_i = 1 \text{ else} \\ Z_i = 0. \end{array} \right]$$

Find sign  $\left( \frac{\epsilon \eta \exp(\pi (X_i - \mu)^2)}{1 - \eta} - 1 \right)$  If sign = 1, above it is true, else false.

2. Suppose we are given values of  $Z_1, \dots, Z_n \in \{0, 1\}$ . Derive an expression for the MLE estimate

$$\operatorname{argmax}_{\mu \in \mathbb{R}} \prod_{i=1}^n \mathbb{P}[X_i | \mu, Z_i]$$

taking negative log-likelihood,

$$\operatorname{argmin}_{\mu \in \mathbb{R}} \sum_{i=1}^n -\ln \mathbb{P}[X_i | \mu, Z_i]$$

$$= \operatorname{argmin}_{\mu \in \mathbb{R}} \underbrace{\sum_{i; Z_i=1} -\ln \epsilon}_{\text{independent of } \mu} + \sum_{i; Z_i=0} \pi (X_i - \mu)^2$$

$$\hat{\mu}_{MLE} = \operatorname{argmin}_{\mu \in \mathbb{R}} \sum_{i; Z_i=0} \pi (X_i - \mu)^2 \quad \left[ \text{let } n = \text{no. of elements whose } Z_i = 0. \right]$$

f.o. optimality w.r.t  $\mu$ , gives

$$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i; Z_i=0} X_i \quad \left[ \hat{\mu}_{MLE} = \frac{1}{n} \sum_{i; Z_i=0} X_i \right] \quad \text{Ans}$$

Note that this allows us to execute alternating optimization to help Melbo solve the problem even in the presence of corruptions. We can initialize  $\mu$  (say randomly), then use part 1 to set  $Z_i$  values for each  $i$  (set  $Z_i = 1$  if  $\mathbb{P}[Z_i = 1 | X_i, \mu] > \mathbb{P}[Z_i = 0 | X_i, \mu]$  else set  $Z_i = 0$ ), then use part 2 to update  $\mu$  given these  $Z_i$  values and then repeat the process till convergence. (5 + 5 = 10 marks)