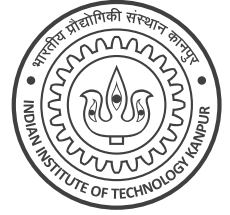


CS 771A: Intro to Machine Learning, IIT Kanpur			Endsem Exam (22 Nov 2022)	
Name				50 marks
Roll No		Dept.		Page 1 of 6

**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>By law of total probability, we have</p> $\begin{aligned} \mathbb{P}[\text{Fall}] &= \mathbb{P}[\text{Fall} \mid \text{Tweet}] \cdot \mathbb{P}[\text{Tweet}] + \mathbb{P}[\text{Fall} \mid \neg \text{Tweet}] \cdot \mathbb{P}[\neg \text{Tweet}] \\ &= \frac{4}{10} \cdot \frac{2}{10} + \frac{1}{10} \cdot \frac{8}{10} = \frac{16}{100} \end{aligned}$		
2	Given three vectors $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^2$ such that $\mathbf{x}^\top \mathbf{z} > \mathbf{x}^\top \mathbf{y}$ , it is always the case that $\ \mathbf{x} - \mathbf{z}\ _2^2 < \ \mathbf{x} - \mathbf{y}\ _2^2$ . Give a proof if True else give a counter example.	F
<p>We can give a counter example even in one dimension. Consider <math>\mathbf{x} = (1, 0) = \mathbf{y}, \mathbf{z} = (100, 0)</math>. We have <math>\mathbf{x}^\top \mathbf{y} = 1 &lt; 100 = \mathbf{x}^\top \mathbf{z}</math>. However, we also have <math>\ \mathbf{x} - \mathbf{z}\ _2^2 = 99^2 &gt; 0 = \ \mathbf{x} - \mathbf{y}\ _2^2</math>.</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 $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ , $\phi(\mathbf{x})^\top \phi(\mathbf{y}) = (1 + \mathbf{x}^\top \mathbf{y})^2$ must use $d \geq 10$ dims. Give a proof if True else give a map using fewer dimensions as a counter example.	F
<p>In general, we do need <math>1 + 2 \cdot 3 + \binom{3}{2} = 10</math> dims. However, for vectors with <math>\pm 1</math> coordinates, the dims encoding <math>x_i^2</math> can be omitted since <math>(+1)^2 = 1 = (-1)^2</math>. The resulting map looks like</p> $\phi(x_1, x_2, x_3) = \sqrt{2}(\sqrt{2}, x_1, x_2, x_3, x_1x_2, x_1x_3, x_2x_3) \in \mathbb{R}^7$ <p>We have</p> $\begin{aligned} \phi(\mathbf{x})^\top \phi(\mathbf{y}) &= 4 + 2x_1y_1 + 2x_2y_2 + 2x_3y_3 + 2x_1x_2y_1y_2 + 2x_1x_3y_1y_3 + 2x_2x_3y_2y_3 \\ &= 1 + x_1^2y_1^2 + x_2^2y_2^2 + x_3^2y_3^2 + 2x_1y_1 + 2x_2y_2 + 2x_3y_3 + 2x_1x_2y_1y_2 + 2x_1x_3y_1y_3 \\ &\quad + 2x_2x_3y_2y_3 \\ &= 1 + (x_1y_1 + x_2y_2 + x_3y_3)^2 + 2(x_1y_1 + x_2y_2 + x_3y_3) = (1 + \mathbf{x}^\top \mathbf{y})^2 \end{aligned}$		

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

Take  $X = \mathbf{1}\mathbf{1}^T = Y$ , both are rank-one, but their sum is  $2 \cdot \mathbf{1}\mathbf{1}^T$  that is rank-one itself.

**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)**

A rich person is less likely to reveal their income than someone with median income. This is because  $\mathbb{P}[Z_i \neq \phi \mid X_i = v]$  is high only if  $|v|$  is small i.e., if  $v$  has large magnitude, never mind its sign, the reveal probability dips.

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)**

A poor person is also less likely to reveal their income than someone with median income. This is because  $\mathbb{P}[Z_i \neq \phi \mid X_i = v]$  is high only if  $|v|$  is small i.e., if  $v$  has large magnitude, never mind its sign, the reveal probability dips.

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)**

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

By law of total probability, we get

$$h(\alpha) \stackrel{\text{def}}{=} \mathbb{P}[Z_i \neq \phi] = \int_{-\infty}^{\infty} \mathbb{P}[Z_i \neq \phi, X_i = v] dv = \frac{1}{\sqrt{1 + \alpha^2}}$$

since

$$\int_{-\infty}^{\infty} \mathbb{P}[Z_i \neq \phi | X_i = v] \cdot \mathbb{P}[X_i = v] dv = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \exp\left(-\frac{(1 + \alpha^2)v^2}{2}\right) dv = \frac{1}{\sqrt{1 + \alpha^2}}$$

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

$$\mathcal{L}(\alpha) = - \sum_{i: Z_i \neq \phi} \ln \mathbb{P}[Z_i \neq \phi, X_i] - \sum_{i: Z_i = \phi} \ln \mathbb{P}[Z_i = \phi]$$

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

By using  $\mathbb{P}[Z_i \neq \phi, X_i] = \mathbb{P}[Z_i \neq \phi | X_i] \cdot \mathbb{P}[X_i]$ , we get

$$\mathcal{L}(\alpha) = - \ln \frac{1}{\sqrt{2\pi}} + \sum_{i: Z_i \neq \phi} \frac{(1 + \alpha^2)X_i^2}{2} - n_0 \cdot \ln \left(1 - \frac{1}{\sqrt{1 + \alpha^2}}\right)$$

where  $n_0 = |i: Z_i = \phi|$ .

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

$$\mathcal{L}'(\alpha) = \alpha \cdot \sum_{i: Z_i \neq \phi} X_i^2 - \frac{n_0}{\sqrt{1 + \alpha^2} - 1} \cdot \frac{\alpha}{1 + \alpha^2}$$

where  $n_0 = |i: Z_i = \phi|$ .

**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)

We have  $f_n(z) = \sum_{x_i \geq z} (x_i - z)$ . Notice that  $f_n(z) \geq 0$  for all  $z$  by definition. However, we see that  $f_n(x_n) = x_n - x_n = 0$ . Thus,  $x_n$  is a minimizer for  $\operatorname{argmin}_{z \in [x_1, x_n]} f_n(z)$ .

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)

Similarly,  $f_0(z) = -\sum_{x_i < z} (x_i - z)$ . Notice that  $f_0(z) \geq 0$  for all  $z$  by definition. However, we see that  $f_0(x_1) = 0$  since it is an empty sum. Thus,  $x_1$  is a minimizer for  $\operatorname{argmin}_{z \in [x_1, x_n]} f_0(z)$ .

5. 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
  - a. We have  $f_k(a) > f_k(b)$  if  $1 \leq j < k$ .
  - b. We have  $f_k(a) < f_k(b)$  if  $k < j < n$ , we have.
  - c. 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)**

CS 771A: Intro to Machine Learning, IIT Kanpur			Endsem Exam (22 Nov 2022)	
Name				50 marks Page 5 of 6
Roll No		Dept.		

Let  $x_j < a < b \leq x_{j+1}$ . Notice that no summation has a zero coefficient when  $k \in [1, n-1]$

$$f_k(a) = \left(\frac{k}{n} - 1\right) \sum_{i=1}^j (x_i - a) + \frac{k}{n} \sum_{i=j+1}^n (x_i - a)$$

$$f_k(b) = \left(\frac{k}{n} - 1\right) \sum_{i=1}^j (x_i - b) + \frac{k}{n} \sum_{i=j+1}^n (x_i - b)$$

giving us  $f_k(a) - f_k(b) = \left[\left(\frac{k}{n} - 1\right)j + \frac{k}{n}(n-j)\right](b-a) = (k-j)(b-a)$ . Since  $b > a$

1. If  $j < k$  i.e., if  $j \leq k-1$ , we get  $f_k(a) > f_k(b)$
2. If  $j > k$  i.e., if  $j \geq k+1$ , we get  $f_k(a) < f_k(b)$
3. If  $j = k$ , we get  $f_k(a) = f_k(b)$

To properly establish the minimizer, some more results are needed (**showing these results is not a part of the problem**). Let  $x_j < c < x_{j+1}$  for some  $j \geq 1$ . We have

$$f_k(x_j) = \left(\frac{k}{n} - 1\right) \sum_{i=1}^{j-1} (x_i - x_j) + \frac{k}{n} \sum_{i=j+1}^n (x_i - x_j)$$

$$f_k(c) = \left(\frac{k}{n} - 1\right) \sum_{i=1}^j (x_i - c) + \frac{k}{n} \sum_{i=j+1}^n (x_i - c)$$

$$\begin{aligned} f_k(x_j) - f_k(c) &= \left[\left(\frac{k}{n} - 1\right)(j-1) + \frac{k}{n}(n-j)\right](c - x_j) - \left(\frac{k}{n} - 1\right)(x_j - c) \\ &= \left[\left(\frac{k}{n} - 1\right)j + \frac{k}{n}(n-j)\right](c - x_j) = (k-j)(c - x_j) \end{aligned}$$

Since  $c > x_j$  by design, we have

1.  $f_k(x_j) > f_k(c)$  whenever  $j < k$
2.  $f_k(x_j) < f_k(c)$  whenever  $j > k$
3.  $f_k(x_k) = f_k(c)$

This finishes the complete argument.

**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  $\mathbb{P}[Z_i = 1 \mid \mu] = \eta$ , and that  $\mathbb{P}[X_i \mid \mu, Z_i = 1] = \epsilon$ , and  $\mathbb{P}[X_i \mid \mu, Z_i = 0] = f(X_i; \mu)$ . 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  $\mathbb{P}[Z_i = 1 \mid X_i, \mu] > \mathbb{P}[Z_i = 0 \mid X_i, \mu]$  or not.

Applying Bayes rule tells us that

$$\mathbb{P}[Z_i = 1 \mid X_i, \mu] \propto \mathbb{P}[X_i \mid \mu, Z_i = 1] \cdot \mathbb{P}[Z_i = 1 \mid \mu] = \epsilon \cdot \eta$$

$$\mathbb{P}[Z_i = 0 \mid X_i, \mu] \propto \mathbb{P}[X_i \mid \mu, Z_i = 0] \cdot \mathbb{P}[Z_i = 0 \mid \mu] = \exp(-\pi(X_i - \mu)^2) \cdot (1 - \eta)$$

Thus,  $\mathbb{P}[Z_i = 1 \mid X_i, \mu] > \mathbb{P}[Z_i = 0 \mid X_i, \mu]$  if  $|X_i - \mu| > \sqrt{\frac{1}{\pi} \ln \left( \frac{(1-\eta)}{\epsilon \cdot \eta} \right)}$

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

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

We have

$$\prod_{i=1}^n \mathbb{P}[X_i \mid \mu, Z_i] = \prod_{i: Z_i=0} \exp(-\pi(X_i - \mu)^2) \cdot \prod_{i: Z_i=1} \epsilon$$

Taking logarithms and setting the derivative w.r.t.  $\mu$  to zero gives us the MLE estimate as

$$\mu_{\text{MLE}} = \frac{1}{|i: Z_i = 0|} \sum_{i: Z_i=0} X_i$$

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