

Name: Roll No.: Dept.: IIT Kanpur  
CS771A (IML)  
Mid-sem Exam

Date: September 23, 2023

**Instructions:****Total: 60 marks**

1. Total duration: **2 hours**. Please write your name, roll number, department on **all pages**.
2. This booklet has 6 pages (5 pages + 1 page for rough work). No part of your answers should be on pages designated for rough work. Additional rough sheets may be provided if needed.
3. Write/mark your answers clearly in the provided space. Please keep your answers precise and concise.
4. Avoid showing very detailed derivations. You may do those on rough sheet and only show key steps.
5. **Answer each question using information provided in the question (no clarifications during the exam). If you want to make any assumptions, please state them in your answer.**

**Section 1** (2 Multiple Choice Questions: Total 4 marks). (Tick/circle all options that you think are true)

1. Which of the following is true about MLE? (1) It is equivalent to finding the best parameters by minimizing the loss function on the training data, (2) It is guaranteed to give the global optimal solution for the parameters, (3) It yields the same answer as maximum-a-posteriori (MAP) estimation if MAP is done using zero-mean Gaussian prior on the parameters, (4) It has the risk of overfitting.
2. Which of the following is true about probabilistic linear regression with Gaussian likelihood and Gaussian prior on the weight vector and assuming all other hyperparameters as fixed? (1) MLE solution for the weight vector is the same as the solution of least squares regression, (2) Posterior distribution of the weight vector is also Gaussian, (3) The mode of the posterior distribution is the same as the MAP solution (4) The mean of the posterior distribution is the same as the MLE solution.

**Section 2** (8 Descriptive Answer Questions: Total 56 marks).

1. You wish to learn a linear regression model with the loss function  $L(\mathbf{w}) = \sum_{i=1}^N |y_i - \mathbf{w}^\top \mathbf{x}_i|$  subject to constraints on  $\mathbf{w} \in \mathbb{R}^D$  that its entries can only be non-negative. Briefly describe an optimization algorithm to learn  $\mathbf{w}$  and write all the necessary mathematical equations/updates required by the algorithm. Your approach must not involve introducing any additional variables to solve this problem. **(6 marks)**.

The loss function is non-differentiable so we will need sub-gradients. As shown in the class, for the absolute loss function, the subgradient w.r.t.  $\mathbf{w}$  at point  $(\mathbf{x}_n, y_n)$  will be  $\mathbf{g}_n = -\text{sign}(y_n - \mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n$  if  $(y_n - \mathbf{w}^\top \mathbf{x}_n) \neq 0$  and  $\mathbf{g}_n = c \mathbf{x}_n$  where  $c \in [-1, +1]$  if  $(y_n - \mathbf{w}^\top \mathbf{x}_n) = 0$ . The overall subgradient will be  $\mathbf{g} = \sum_{n=1}^N \mathbf{g}_n$ . Given these subgradient, we can perform subgradient descent where in each iteration, we update  $\mathbf{w}$  as  $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}_t$  where  $\mathbf{g}$  denotes the (sub)gradient in iteration  $t$ .

To handle the non-negativity constraint, we can use projected gradient descent which in this case reduces to simply setting the negative entries of  $\mathbf{w}$  to 0 in every iteration of the subgradient descent algorithm.

2. Suppose we are given  $N$  training examples  $\{(x_i, y_i)\}_{i=1}^N$  with each input  $x_n \in \mathbf{R}$  being a scalar value and each output  $y_n \in \{1, 2, \dots, K\}$ . Assume the inputs from class  $c$  are generated i.i.d. by univariate Gaussian  $\mathcal{N}(x|\mu_c, \sigma_c^2)$  and, for simplicity, assume that only  $\mu_c$  is unknown and  $\sigma_c^2$  is known.

- Derive (only show key steps) and write down the expression for the MLE solution of  $\mu_c$ . **(3 marks)**  
Since the inputs are assumed generated i.i.d., the log likelihood will be the sum of the log likelihoods of the inputs from the  $K$  classes:

$$\text{LL} = \sum_{c=1}^K \sum_{i: y_i=c} \log \mathcal{N}(x_i | \mu_c, \sigma_c^2) = \sum_{c=1}^K \sum_{i: y_i=c} \log \frac{1}{\sqrt{2\pi\sigma_c^2}} \exp\left(-\frac{(x_i - \mu_c)^2}{2\sigma_c^2}\right)$$

Ignoring  $\sigma_c$  and other constants, the negative log-lik. is  $\text{NLL} = \sum_{c=1}^K \sum_{i: y_i=c} \frac{(x_i - \mu_c)^2}{2\sigma_c^2}$ . Taking the NLL's partial derivative w.r.t.  $\mu_c$  and using first-order optimality, the MLE solution  $\hat{\mu}_c = \frac{\sum_{i: y_i=c} x_i}{N_c}$  which is just the mean of inputs from class  $c$ .

Name: Roll No.: Dept.: 

- Given the MLE solution of  $\{\mu_c\}_{c=1}^K$ , to predict the label  $y_*$  for a new test input  $x_*$ , suppose this model uses the following rule at test time: predict  $y_*$  to be the class under which the test input  $x_*$  has the largest probability density. Compare and contrast this model with the Learning with the Prototypes (LwP) model. What would be the advantage of this model as compared to LwP? (5 marks)

From the given prediction rule,  $y_* = \operatorname{argmax}_{c=1}^K \mathcal{N}(x_* | \mu_c, \sigma_c) = \operatorname{argmax}_{c=1}^K \frac{1}{\sqrt{2\pi\sigma_c^2}} \exp(-\frac{(x_* - \mu_c)^2}{2\sigma_c^2})$ . Taking log (since log is monotonic), this is equivalent to  $y_* = \operatorname{argmax}_{c=1}^K -\frac{(x_* - \mu_c)^2}{2\sigma_c^2} = \operatorname{argmin}_{c=1}^K \frac{(x_* - \mu_c)^2}{2\sigma_c^2}$ .

So  $y_*$  is the class for which  $\frac{(x_* - \mu_c)^2}{2\sigma_c^2}$  is the smallest. Note that this is the squared Euclidean distance of  $x_*$  from the mean of class  $c$ , dividing by the variance of class  $c$ . If variance of each class is equal, then this prediction rule is exactly the same as LwP. Otherwise, this prediction rule depends on both the mean  $\mu_c$  as well as the variance (spread)  $\sigma_c^2$  of class  $c$ .

As per this prediction rule, in addition to a small Euclidean distance  $(x_* - \mu_c)^2$ , if the variance  $\sigma_c^2$  of class  $c$  is also larger (i.e., this class has a large spread), the probability of  $x_*$  belonging to class  $c$  will be even higher.

- Given a set  $P$  with  $N$  inputs  $\{\mathbf{x}_i\}_{i=1}^N$  and another set  $Q$  with  $M$  inputs  $\{\mathbf{z}_j\}_{j=1}^M$ , one notion of distance between these two sets is the squared Euclidean distance  $\|\phi_P - \phi_Q\|^2$  between their means  $\phi_P$  and  $\phi_Q$  in a feature space  $\phi$  defined by a kernel function  $k$ . Derive the expression of this distance. Your final answer must be written only in terms of  $k$  and not in terms of the feature mapping  $\phi$ . (6 marks)

$\phi_P = \frac{\sum_{i=1}^N \phi(\mathbf{x}_i)}{N}$  and  $\phi_Q = \frac{\sum_{j=1}^M \phi(\mathbf{z}_j)}{M}$ . Thus  $\|\phi_P - \phi_Q\|^2 = \|\phi_P\|^2 + \|\phi_Q\|^2 - 2\phi_P^\top \phi_Q$ , where

$$\|\phi_P\|^2 = \phi_P^\top \phi_P = \frac{1}{N^2} \sum_{\mathbf{x}_i, \mathbf{x}_j \in P} \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) = \frac{1}{N^2} \sum_{\mathbf{x}_i, \mathbf{x}_j \in P} k(\mathbf{x}_i, \mathbf{x}_j)$$

$$\|\phi_Q\|^2 = \phi_Q^\top \phi_Q = \frac{1}{M^2} \sum_{\mathbf{z}_i, \mathbf{z}_j \in Q} \phi(\mathbf{z}_i)^\top \phi(\mathbf{z}_j) = \frac{1}{M^2} \sum_{\mathbf{z}_i, \mathbf{z}_j \in Q} k(\mathbf{z}_i, \mathbf{z}_j)$$

$$\phi_P^\top \phi_Q = \frac{1}{NM} \sum_{\mathbf{x}_i \in P, \mathbf{z}_j \in Q} \phi(\mathbf{x}_i)^\top \phi(\mathbf{z}_j) = \frac{1}{NM} \sum_{\mathbf{x}_i \in P, \mathbf{z}_j \in Q} k(\mathbf{x}_i, \mathbf{z}_j)$$

- Suppose we have rolled a  $K$  faced die a total of  $N$  times and the number of times each face showed is given by  $N_1, N_2, \dots, N_K$ . Assume each outcome  $\mathbf{x}_n$  (a one-hot vector of length  $K$ ) of the die roll to be a random variable modeled by a multinoulli  $p(\mathbf{x}_n | \boldsymbol{\pi})$  where  $\boldsymbol{\pi} = \{\pi_1, \pi_2, \dots, \pi_K\}$  is the parameter of the multinoulli distribution. Note that the multinomial distribution also has the constraint  $\sum_{i=1}^K \pi_i = 1$ .

Using Lagrange multiplier based constrained optimization, derive the MLE solution of  $\{\pi_i\}_{i=1}^K$ . (8 marks)  
For the  $N$  observations drawn i.i.d. from multinoulli with parameters  $\boldsymbol{\pi}$ , the log likelihood

$$\mathcal{L} = \log \prod_{n=1}^N \prod_{i=1}^K \pi_i^{x_{ni}} = \sum_{n=1}^N \sum_{i=1}^K x_{ni} \log \pi_i$$

We also have the equality constraint  $1 - \sum_{i=1}^K \pi_i = 0$ . Using a Lagrange multiplier  $\lambda$ , the Lagrangian for this problem will be

$$\mathcal{Q} = \log \prod_{n=1}^N \prod_{i=1}^K \pi_i^{x_{ni}} = \sum_{n=1}^N \sum_{i=1}^K x_{ni} \log \pi_i + \lambda(1 - \sum_{i=1}^K \pi_i)$$

Taking partial derivative w.r.t.  $\pi_i$  and setting it to zero gives  $\frac{1}{\pi_i} \sum_{n=1}^N x_{ni} - \lambda = 0$ . Thus  $\pi_i = \frac{1}{\lambda} \sum_{n=1}^N x_{ni}$ . Now what is the value of  $\lambda$ ? To get  $\lambda$ , summing both sides of  $\lambda \pi_i = \sum_{n=1}^N x_{ni}$  over all  $i$ , and using the fact  $\sum_{i=1}^K \pi_i = 1$ , we can see that  $\lambda = N$ .

Note: We don't need to use the additional non-negativity constraints on each  $\pi_i$  because the  $\log \pi_i$  term in the objective already ensures that.

Name: Roll No.: Dept.: 

5. Define an  $\epsilon$ -insensitive loss for linear regression as  $\ell_\epsilon(y_n, \hat{y}_n) = 0$  if  $|y_n - \hat{y}_n| \leq \epsilon$ , otherwise  $\ell_\epsilon(y_n, \hat{y}_n) = (|y_n - \hat{y}_n| - \epsilon)^2$ . Here  $\hat{y}_n = \mathbf{w}^\top \mathbf{x}_n$  denotes the model's prediction and  $y_n$  is the true label. Write down the expression for likelihood  $p(y_n | \mathbf{x}_n, \mathbf{w})$  that corresponds to this loss function, and plot its PDF. (6 marks)
- We know that the negative log-likelihood corresponds to the loss function. Using this fact, the likelihood is (up to a proportionality constant) the exponential of the negative of loss. Thus  $p(y_n | \mathbf{w}, \mathbf{x}_n) = C \exp(-\ell_\epsilon(y_n, \hat{y}_n))$ . For the two cases, we have

$$p(y_n | \mathbf{w}, \mathbf{x}_n) = \begin{cases} C & \text{if } |y_n - \hat{y}_n| \leq \epsilon \\ C \exp(-(|y_n - \hat{y}_n| - \epsilon)^2), & \text{otherwise} \end{cases}$$

Using the fact that the area under the curve is 1, and the curve is symmetric around  $|y_n - \hat{y}_n| = 0$ , we have  $C \times 2\epsilon + 2C \int_\epsilon^\infty \exp(-(|y_n - \hat{y}_n| - \epsilon)^2) d|y_n - \hat{y}_n| = 1$ . Assuming the integral equals  $K$ ,  $C = 1/(2\epsilon + K)$

6. The ridge regression problem has the solution  $\mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$ , where  $\mathbf{X}$  is  $N \times D$  feature matrix (row  $i$  contains the input  $\mathbf{x}_i$ ) and  $\mathbf{y}$  is the  $N \times 1$  response vector, and  $\lambda > 0$  is the regularization hyperparameter. Using the result  $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1}$  (due to Woodbury matrix identity), show that the ridge regression model can be kernelized. More specifically, show that, given a test input  $\mathbf{x}_*$ , we can compute a kernelized prediction which is given by  $y_* = \sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x}_*)$ . (6 marks)
- The ridge regression prediction can be written as

$$y_* = \mathbf{w}^\top \mathbf{x}_* = \mathbf{x}_*^\top \mathbf{w} = \mathbf{x}_*^\top \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{y} = [\mathbf{x}_*^\top \mathbf{x}_1, \mathbf{x}_*^\top \mathbf{x}_2, \dots, \mathbf{x}_*^\top \mathbf{x}_N] (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

Note that  $\mathbf{x}_*^\top \mathbf{X}^\top = [\mathbf{x}_*^\top \mathbf{x}_1, \mathbf{x}_*^\top \mathbf{x}_2, \dots, \mathbf{x}_*^\top \mathbf{x}_N]$  is an  $1 \times N$  vector with its  $i$ -th entry being the simple dot product similarity of  $\mathbf{x}_*$  with the  $i$ -th training example  $\mathbf{x}_i$ . In the kernelized version, this will be replaced by  $[\phi(\mathbf{x}_*)^\top \phi(\mathbf{x}_1), \phi(\mathbf{x}_*)^\top \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_*)^\top \phi(\mathbf{x}_N)] = [k(\mathbf{x}_*, \mathbf{x}_1), k(\mathbf{x}_*, \mathbf{x}_2), \dots, k(\mathbf{x}_*, \mathbf{x}_N)] = \mathbf{v}_*$  where  $\mathbf{v}_*$  is the  $1 \times N$  vector with the  $i$ -th entry being equal to  $k(\mathbf{x}_*, \mathbf{x}_i)$ .

Also note that  $\mathbf{X} \mathbf{X}^\top$  is the  $N \times N$  matrix with its  $(i, j)$ -th entry being equal to  $\mathbf{x}_i^\top \mathbf{x}_j$ . In the kernelized version, we can replace  $(\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$  with  $(\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y} = \boldsymbol{\alpha}$ , which is an  $N \times 1$  vector and  $K_{ij} = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j)$ .

Combining these two, we have  $y_* = \mathbf{v}_* \boldsymbol{\alpha} = \boldsymbol{\alpha}^\top \mathbf{v}_* = \sum_{i=1}^N \alpha_i k(\mathbf{x}_*, \mathbf{x}_i)$ .

7. Consider a logistic regression model  $p(y_n | \mathbf{x}_n, \mathbf{w}) = \frac{1}{1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}$ , with a zero-mean Gaussian prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \lambda^{-1} \mathbf{I})$ . Note that this loss function for logistic regression assumes  $y_n \in \{-1, +1\}$  instead of  $\{0, 1\}$ . Show that the MAP estimate for  $\mathbf{w}$  can be written as  $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$  where each  $\alpha_n$  itself is a function of  $\mathbf{w}$ . Based on the expression of  $\alpha_n$ , you would see that it has a precise meaning. Briefly explain what  $\alpha_n$  means, and also briefly explain why the result  $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$  makes sense. (8 marks)

The objective function (log-posterior = log-likelihood + log-prior) is  $\sum_{n=1}^N \log p(y_n | \mathbf{x}_n, \mathbf{w}) + \log p(\mathbf{w}) = -\sum_{n=1}^N \log(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)) - \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$ . Taking its derivative of w.r.t.  $\mathbf{w}$  and setting it to zero gives

$$\mathbf{w} = \frac{1}{\lambda} \sum_{n=1}^N \frac{1}{1 + \exp(y_n \mathbf{w}^\top \mathbf{x}_n)} y_n \mathbf{x}_n = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

In the above expression, we can think of  $\alpha_n = \frac{1}{1 + \exp(y_n \mathbf{w}^\top \mathbf{x}_n)} = 1 - \frac{1}{1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n)}$  as denoting the importance of the  $n$ -th training example in the final solution  $\mathbf{w}$ . Note that this term is also the probability of **mis-classification**. Thus the training examples that are harder (thus having a higher probability of being mis-classified) are deemed as more important and contribute more to the solution of  $\mathbf{w}$ . Note that solving for  $\mathbf{w}$  in logistic regression requires an iterative optimization (no closed form solution), and in each iteration, these  $\alpha_n$ 's determine how important the  $n$ -th training example is for the current optimal  $\mathbf{w}$ .

Name: Roll No.:  Dept.: 

8. Consider the dual problem for soft-margin linear SVM:  $\arg \max_{0 \leq \alpha \leq C} f(\alpha)$ , where  $f(\alpha) = \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top \mathbf{G} \alpha$ ,  $\mathbf{G}$  is an  $N \times N$  matrix such that  $G_{nm} = y_n y_m \mathbf{x}_n^\top \mathbf{x}_m$  (assume labels as -1/+1), and  $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_N]$  are the Lagrange multipliers. Given the optimal  $\alpha$ , the SVM weight vector is  $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$ .

Your goal is to derive a **co-ordinate ascent** procedure for the vector  $\alpha$ , such that each iteration updates a uniformly randomly chosen entry  $\alpha_n$  of the vector  $\alpha$ . However, instead of updating  $\alpha$  via standard co-ordinate descent as  $\alpha_n = \alpha_n + \eta g_n$  where  $g_n$  denotes the  $n$ -th entry of the gradient vector  $\nabla_{\alpha} f(\alpha)$ , we will update it as  $\alpha_n = \alpha_n + \delta_*$  where  $\delta_* = \arg \max_{\delta} f(\alpha + \delta \mathbf{e}_n)$  and  $\mathbf{e}_n$  denotes a vector of all zeros except a 1 at entry  $n$ . Essentially, this will give the new  $\alpha_n$  that guarantees the maximum increase in  $f$ , with all other  $\alpha_n$ 's fixed at their current value. Derive the expression for  $\delta_*$ .

Note that your expression for  $\delta_*$  should be such that the constraint  $0 \leq \alpha_n \leq C$  is maintained. (8 marks)

$$\delta_* = \arg \max_{\delta} f(\alpha + \delta \mathbf{e}_n) = \arg \max_{\delta} (\alpha + \delta \mathbf{e}_n)^\top \mathbf{1} - \frac{1}{2} (\alpha + \delta \mathbf{e}_n)^\top \mathbf{G} (\alpha + \delta \mathbf{e}_n)$$

Expanding the above, the problem becomes

$$\delta_* = \arg \max_{\delta} (\alpha_n + \delta) + \sum_{m \neq n} \alpha_m - \frac{1}{2} \alpha^\top \mathbf{G} \alpha - \alpha^\top \mathbf{G} \mathbf{e}_n \delta - \frac{1}{2} G_{nn} \delta^2$$

The above is a quadratic function of  $\delta$  with a unique maxima. Taking the first derivative w.r.t.  $\delta$  and setting it to zero gives  $1 - (\mathbf{G} \alpha)_n - G_{nn} \delta = 0$  where  $(\mathbf{G} \alpha)_n$  denotes the  $n$ -th entry of the vector  $\mathbf{G} \alpha$ . This yields  $\delta_* = \frac{1 - (\mathbf{G} \alpha)_n}{G_{nn}}$ .

The update  $\alpha_n = \alpha_n + \delta_*$  must also satisfy the constraint  $0 \leq \alpha_n + \delta_* \leq C$ , i.e.,  $-\alpha_n \leq \delta_* \leq C - \alpha_n$ . We can use projection operator (used in projected gradient descent) on  $\delta_*$  to ensure this constraint. Whenever  $\delta_* = \frac{1 - (\mathbf{G} \alpha)_n}{G_{nn}}$  becomes smaller than  $-\alpha_n$ , we set it equal to  $-\alpha_n$  and whenever  $\delta_* = \frac{1 - (\mathbf{G} \alpha)_n}{G_{nn}}$  becomes larger than  $C - \alpha_n$ , we set it equal to  $C - \alpha_n$ . Mathematically, this can be written as an equation  $\delta_* = \max \left[ -\alpha_n, \min(C - \alpha_n, \frac{1 - (\mathbf{G} \alpha)_n}{G_{nn}}) \right]$ .

### Some distributions and their properties:

- For  $x \in \mathbb{R}$ , the PDF of univariate Gaussian:  $\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\{-\frac{(x-\mu)^2}{2\sigma^2}\}$ . If using precision  $\beta = 1/\sigma^2$ , the PDF is  $\mathcal{N}(x|\mu, \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp\{-\frac{\beta}{2}(x - \mu)^2\}$ .
- Given  $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$ , s.t.  $\sum_{k=1}^K \pi_k = 1$ , the multinoulli is defined as  $\text{multinoulli}(\mathbf{x}|\boldsymbol{\pi}) = \prod_{i=1}^K \pi_i^{x_i}$ , where the random variable  $\mathbf{x} = [x_1, x_2, \dots, x_K]$  is a one-hot vector.

### Some other useful results:

- $\frac{\partial}{\partial \mathbf{x}} \mathbf{a}^\top \mathbf{x} = \mathbf{a}$ ,  $\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = 2\mathbf{A} \mathbf{x}$ .  $|x| = \max\{x, -x\}$ . If a function  $h(x) = \max\{f(x), g(x)\}$  then, at any point  $x_*$ , (1) if  $f(x_*) > g(x_*)$  then the subgradient  $\partial h(x_*) = \partial f(x_*)$ , (2) if  $g(x_*) > f(x_*)$  then  $\partial h(x_*) = \partial g(x_*)$ , and (3) if  $f(x_*) = g(x_*)$  then  $\partial h(x_*) = \{\alpha a + (1 - \alpha)b\}$ , where  $a \in \partial f(x_*)$ ,  $b \in \partial g(x_*)$ , and  $\alpha \in [0, 1]$ .

Name:

Roll No.:

Dept.:

**IIT Kanpur**  
**CS771A (IML)**  
**Mid-sem Exam**

*Date:* September 23, 2023

---

FOR ROUGH WORK ONLY