### CSE-4819/5820 Introduction to Machine Learning
### Written Assignment 1
### Linear Algebra; Convex Optimization; Linear Regression
Total: 50pt

Due: End of Day (23:59), September 20th

**Requirements:**

1. The written assignment will be graded based on correctness, accuracy and clarity.
2. Please prepare your answers, and submit the final assignment in a **pdf** file through HuskyCT assignment portal. It is recommended that you prepare the assignment using MS Word and then convert it into pdf. Please put your **full name** at the beginning of the document for ease of grade recording.
3. Please explicitly indicate the ids of the questions you are answering for ease of grading.
4. Even if your answer is not correct, you may still get certain partial marks based on your calculation/analysis process. Please present the necessary calculation process if any.
5. Each student will have a total of 5 free late (calendar) days to use for homework. Each 24 hours or part thereof that an assignment is late uses up one full late day. Please note: once these late days are exhausted, no late assignments will be accepted for any reason. Students are highly encouraged to reserve your late days for unavoidable emergencies, planned travel, etc.

**Q1. Linear Algebra and Probability** (22pt)

(1) (8pt) Given the following four vectors,

$x_1 = [0, 0.2, 1.0, 2.2]$

$x_2 = [0.7, 0.2, 0.5, 2.0]$

$x_3 = [0, 1.0, 1.5, 2.2]$

$x_4 = [0.8, 0.1, 1.2, 2.0]$

Which point is closest to $x_1$ under each of the following norms?

a) $L_0$

b) $L_1$

c) $L_2$

d) $L_\infty$

Which point is closest to $x_1$ under each of the following norms?

a) $L_0$   $x_3$ with distance $= 2$

b) $L_1$   $x_3$ and $x_4$ with distance $= 1.3$

c) $L_2$   $x_4$ with distance $= 0.85$

d) $L_\infty$   $x_2$ with distance $= 0.7$

(2) (4pt) If $X \sim N(\mu, \sigma^2)$, $E[X] = \mu$, Var[X] $=\sigma^2$, and $E[X^2] = \mu^2 + \sigma^2$. Also, recall that expectation is linear, so it obeys the following three properties:

$E[X + c] = E[X] + c$ for any constant c,

$E[X + Y] = E[X] + E[Y]$,

$E[aX] = aE[X]$ for any constant a.

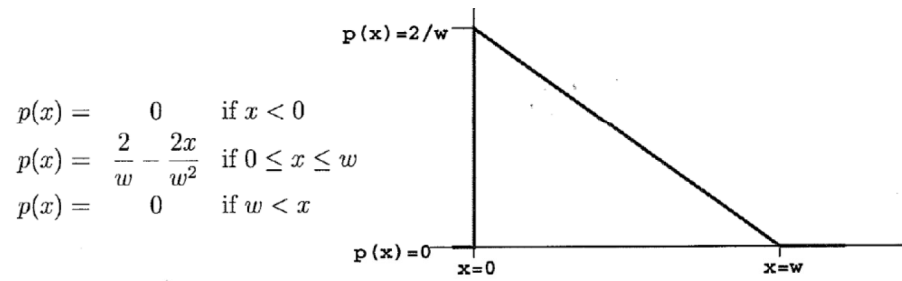We note that if $X$ and $X'$ are independent, then $E[XX'] = E[X]E[X']$.

Consider two points (sampled independently) from the same class follow: $X \sim N(\mu_1, \sigma^2)$ and $X' \sim N(\mu_1, \sigma^2)$.

What is the expected squared distance between them, i.e., $E[(X - X')^2]$?

What is the expected squared distance between them, i.e., $E[(X - X')^2]$?

$$E[(X-X')^2] = E[X^2 - 2XX' + X'^2] = E[X^2] - 2E[XX'] + E[X'^2]$$
$$= \mu_1^2 + \sigma^2 - 2E[X]E[X'] + \mu_1^2 + \sigma^2 = 2\sigma^2$$

(3) (4pt) Consider the probability density function shown in the following figure and equations.

$$p(x) = 2/w$$

$$p(x) = \begin{cases} 0 & \text{if } x < 0 \\ \dfrac{2}{w} - \dfrac{2x}{w^2} & \text{if } 0 \leq x \leq w \\ 0 & \text{if } w < x \end{cases}$$

$$p(x) = 0$$

$x=0$      $x=w$

(2pt) Which _one_ of the following expressions is true?

(a) $E[X] = \int_{x=-\infty}^{\infty} w(\frac{2}{w} - \frac{2x}{w^2})dx$

(b) $E[X] = \int_{x=0}^{w} x(\frac{2}{w} - \frac{2x}{w^2})dx$

(c) $E[X] = \int_{x=0}^{w} w(\frac{2}{w} - \frac{2x}{w^2})dx$

(d) $E[X] = \int_{x=-\infty}^{\infty} (\frac{2}{w} - \frac{2x}{w^2})dx$

(2pt) What is $p(x = 1 | w = 2)$? 0.5

(4) (6pt) Consider a feature $x$ which is a continuous random variable with possible outcomes being all the nonnegative real numbers. The random variable follows a distribution with the following probability density function (PDF):

$$p(x \mid \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & x < 0, \end{cases}$$

where the parameter $\lambda$ of the distribution is a positive real number.

Given a data set $X = \{x_1, x_2, \ldots, x_N\}$ drawn *i.i.d.* (independent and identically distributed) from the distribution, derive the maximum likelihood estimate (MLE) of $\lambda$ based on $X$.

The log likelihood of $\lambda$ given $X$ is

$$\mathcal{L}(\lambda \mid X) = N \ln \lambda - \lambda \sum_{\ell=1}^{N} x_\ell$$

To maximize $\mathcal{L}$, we set its derivative with respect to $\lambda$ to zero,

$$\frac{d\mathcal{L}}{d\lambda} = \frac{N}{\lambda} - \sum_{\ell=1}^{N} x_\ell = 0$$

Then

$$\hat{\lambda} = \frac{1}{\frac{1}{N} \sum_{\ell=1}^{N} x_\ell}$$

**Q2. Introduction to Optimization** (16pt)

(1) (10pt) Please justify if the following statement is correct or not (2 pt for correct T/F; 3 pt each for correct explanation).

   (a) In machine learning, the optimization problem we are solving to train a model is always a maximization problem.

Q2 (a) True or False. True: ☐; Any optimization with minimization can be converted int a maximization problem. False: Sometimes we are maximizing a likelihood, and (by flipping the sign) sometimes we are minimize a cost function.

   (b) For any constant $c \in R$, $f(x) = \frac{x^2}{c-x}$ is convex on $-\infty < x < c$.

Q2(b)   $\frac{x^2}{c-x} = -\frac{x^2}{x-c} = -\left[x+c+\frac{c^2}{x-c}\right]$

$= -(x+c) - \frac{c^2}{x-c}$

Each of $-(x+c)$ and $-\frac{c^2}{x-c}$ is a convex function

on $-\infty < x < c$.

(2) (6pt) Use the method of Lagrange multipliers to find the maximum values of the objective function. Please provide the maximum values and the corresponding variables.

objective function: maximize $f(x, y) = 6xy$

subject to: $x^2/9 + y^2/16 = 1$.

(2) $L = f(x,y) + \lambda\left(\dfrac{x^2}{9} + \dfrac{y^2}{16} - 1\right)$

$\dfrac{\partial L}{\partial x} = 6y + \dfrac{2\lambda x}{9} = 0 \qquad \dfrac{\partial L}{\partial y} = 6x + \dfrac{2\lambda y}{16} = 0, \qquad \dfrac{\partial L}{\partial \lambda} = \dfrac{x^2}{9} + \dfrac{y^2}{16} - 1 = 0$

$\Rightarrow \lambda = \dfrac{-27y}{x} = \dfrac{-48x}{y} \qquad \Rightarrow \dfrac{x^2}{9} = \dfrac{y^2}{16}$

Since $\dfrac{x^2}{9} + \dfrac{y^2}{16} = 1 \qquad \Rightarrow \quad x = \pm\dfrac{3\sqrt{2}}{2} \quad y = \pm 2\sqrt{2}$

For maximizing $6xy$,

$(x, y)$ is either $\left(\dfrac{3\sqrt{2}}{2}, 2\sqrt{2}\right)$ or $\left(-\dfrac{3\sqrt{2}}{2}, -2\sqrt{2}\right)$

maximum value of $6xy$ is $36$.

**Q3. Linear Regression** (12pt)

Given known $\mathbf{X} \in R^{n\times d}$, $\mathbf{y} \in R^{n\times 1}$, and unknown $\mathbf{w} \in R^{d\times 1}$, $\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon$, where $\epsilon \sim N(0, \sigma^2 I)$.
The task is to estimate $\mathbf{w}$.

(a) (6pt) Please write down the loss function for the linear regression. Then derive the closed form estimation for $\mathbf{w}$ based on the least square method. Note that the derive process is required and we assume that $\mathbf{X}^T\mathbf{X}$ is invertible, $i.e.$, $(\mathbf{X}^T\mathbf{X})^{-1}$ exists.

(b) (6pt) Given $\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 2 & 2 \\ 1 & 3 \end{bmatrix}$, and $\mathbf{y} = \begin{bmatrix} 5 \\ 3 \\ 2 \end{bmatrix}$. Using the closed form estimation for $\mathbf{w}$ based on

the least square method, please compute $\mathbf{X}^T\mathbf{X}$, $\mathbf{X}^T\mathbf{y}$ and the estimated $\mathbf{w}$.

$$L(w) = \|y - Xw\|_2^2$$

or $\displaystyle\sum_{i=1}^{n} (y_i - w^T x_i)^2$

or $L(w) = (y - Xw)^T (y - Xw)$

$$L(w) = (y - Xw)^T (y - Xw)$$
$$= (y^T - w^T X^T)(y - Xw)$$
$$= y^T y - y^T Xw - w^T X^T y + w^T X^T Xw$$

$$\frac{\partial L(w)}{\partial w} = 0 - 2X^T y + 2X^T Xw = 0.$$

$$\Rightarrow X^T Xw = X^T y \qquad \Rightarrow w = (X^T X)^{-1} X^T y$$

$$X^T X = \begin{bmatrix} 1 & 2 & 1 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 2 & 2 \\ 1 & 3 \end{bmatrix} = \begin{bmatrix} 6 & 8 \\ 8 & 14 \end{bmatrix}$$

$$X^T y = \begin{bmatrix} 1 & 2 & 1 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 5 \\ 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 13 \\ 17 \end{bmatrix}$$

$$\hat{w} = (X^T X)^{-1} X^T y = \begin{bmatrix} 2.3 \\ -0.1 \end{bmatrix}$$