

FIRST NAME: _____ LAST NAME: _____

STUDENT NUMBER: _____

**ECE 421F — Introduction to Machine Learning
MidTerm Examination**

**Wed Oct 16th, 2019
4:10-6:00 p.m.**

Instructor: Ashish Khisti

Circle your tutorial section:

- TUT0101 Thu 1-3
- TUT0102 Thu 4-6

Instructions

- Please read the following instructions carefully.
- You have 1 hour fifty minutes (1:50) to complete the exam.
- Please make sure that you have a complete exam booklet.
- Please answer *all* questions. Read each question carefully.
- The value of each question is indicated. Allocate your time wisely!
- No additional pages will be collected beyond this answer book. You may use the reverse side of each page if needed to show additional work.
- This examination is closed-book; One 8.5×11 aid-sheet is permitted. A non-programmable calculator is also allowed.
- Good luck!

1. (40 MARKS) Consider a multi-class linear classification problem where the data points are two dimensional, i.e., $\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$ and the labels $y \in \{1, 2, 3\}$. Throughout this problem consider the data-set with following five points:

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3), (\mathbf{x}_4, y_4), (\mathbf{x}_5, y_5)\}$$

where the input data-vectors are given by:

$$\mathbf{x}_1 = (-1, 0)^T, \quad \mathbf{x}_2 = (1, 0)^T, \quad \mathbf{x}_3 = (1, 1)^T, \quad \mathbf{x}_4 = (-1, 1)^T, \quad \mathbf{x}_5 = (0, 3)^T$$

and the associated labels are given by

$$y_1 = 1, \quad y_2 = 2, \quad y_3 = 2, \quad y_4 = 1, \quad y_5 = 3$$

Our aim is to find a linear classification rule that classifies this dataset.

10 marks

- (a) Suppose we implement the perceptron learning algorithm for binary classification that finds a perfect classifier separating the data points between the two sets: $\mathcal{S}_1 = \{(\mathbf{x}_1, y_1), (\mathbf{x}_4, y_4)\}$ and $\mathcal{S}_2 = \{(\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3)\}$.

Assume that the initial weight vector $\mathbf{w} = (0, 0, 0)^T$, that each point that falls on the boundary is treated as a mis-classified point and the algorithm visits the points in the following order:

$$\mathbf{x}_1 \rightarrow \mathbf{x}_2 \rightarrow \mathbf{x}_3 \rightarrow \mathbf{x}_4 \rightarrow \mathbf{x}_1 \rightarrow \mathbf{x}_2 \dots$$

until it terminates. Show the output of the perceptron algorithm in each step and sketch the final decision boundary when the algorithm terminates. [Important: When applying the perceptron update, recall that you have to transform the data vectors to include the constant term i.e., $\mathbf{x}_1 = (-1, 0)^T$ must be transformed to $\tilde{\mathbf{x}}_1 = (1, -1, 0)^T$ etc.]

$$\begin{aligned} \mathbf{x}_1, \mathbf{x}_4 &\rightarrow y = 1 \text{ (+ve class)} \\ \mathbf{x}_2, \mathbf{x}_3 &\rightarrow y = -1 \text{ (-ve class)} \end{aligned}$$

$$\text{Sign}(x) = \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{if } x < 0 \\ 0 & \text{if } x = 0 \end{cases}$$

Step 1: $\hat{y}_1 = \text{Sign}(\mathbf{w}^{(0)T} \tilde{\mathbf{x}}_1)$

$$\begin{aligned} &= \text{Sign}\left([0 \ 0 \ 0] \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}\right) \\ &= \text{Sign}(0) \\ &= 0 \rightarrow \text{misclassified} \end{aligned}$$

$$\begin{aligned} \mathbf{w}^{(1)} &= \mathbf{w}^{(0)} + (y_1 - \hat{y}_1) \tilde{\mathbf{x}}_1 \\ &= \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + 1 \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \end{aligned}$$

Step 2: $\hat{y}_2 = \text{Sign}(\mathbf{w}^{(1)T} \tilde{\mathbf{x}}_2)$

$$\begin{aligned} &= \text{Sign}\left([1 \ -1 \ 0] \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}\right) \\ &= \text{Sign}(0) = 0 \rightarrow \text{misclassified} \end{aligned}$$

$$\begin{aligned} \mathbf{w}^{(2)} &= \mathbf{w}^{(1)} + y_2 \tilde{\mathbf{x}}_2 \\ &= \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} + (-1) \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ -2 \\ 0 \end{bmatrix} ; \end{aligned}$$

[continue part (a) here]

$$\begin{aligned}
 \text{Step 3: } \hat{y}_3 &= \text{sign}(\bar{w}^{(2)T} \hat{x}_3) \\
 &= \text{sign}\left([0 \ -2 \ 0] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}\right) \\
 &= \text{sign}(0 - 2 + 0) \\
 &= -1 \rightarrow \text{classified } \bar{w}^3 = \bar{w}^2
 \end{aligned}$$

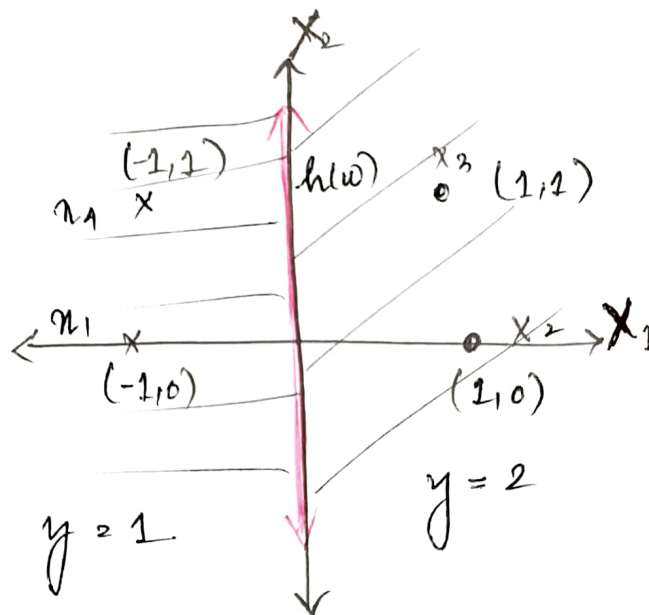
$$\begin{aligned}
 \text{Step 4: } \hat{y}_4 &= \text{sign}(\bar{w}^{(3)T} \hat{x}_4) \\
 &= \text{sign}\left([0 \ -2 \ 0] \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}\right) \\
 &= 1 \rightarrow \text{classified } \bar{w}^4 = \bar{w}^3
 \end{aligned}$$

$$\begin{aligned}
 \text{Step 5: } \hat{y}_1 &= \text{sign}(\bar{w}^{(4)T} \hat{x}_1) \\
 &= \text{sign}\left([0 \ -2 \ 0] \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}\right) \\
 &= +1 \rightarrow \text{classified } \bar{w}^5 = \bar{w}^4
 \end{aligned}$$

$$\begin{aligned}
 \text{Step 6: } \hat{y}_2 &= \text{sign}(\bar{w}^{(5)T} \hat{x}_2) \\
 &= \text{sign}\left([0 \ -2 \ 0] \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}\right) \\
 &= -1 \rightarrow \text{classified } \bar{w}^6 = \bar{w}^5
 \end{aligned}$$

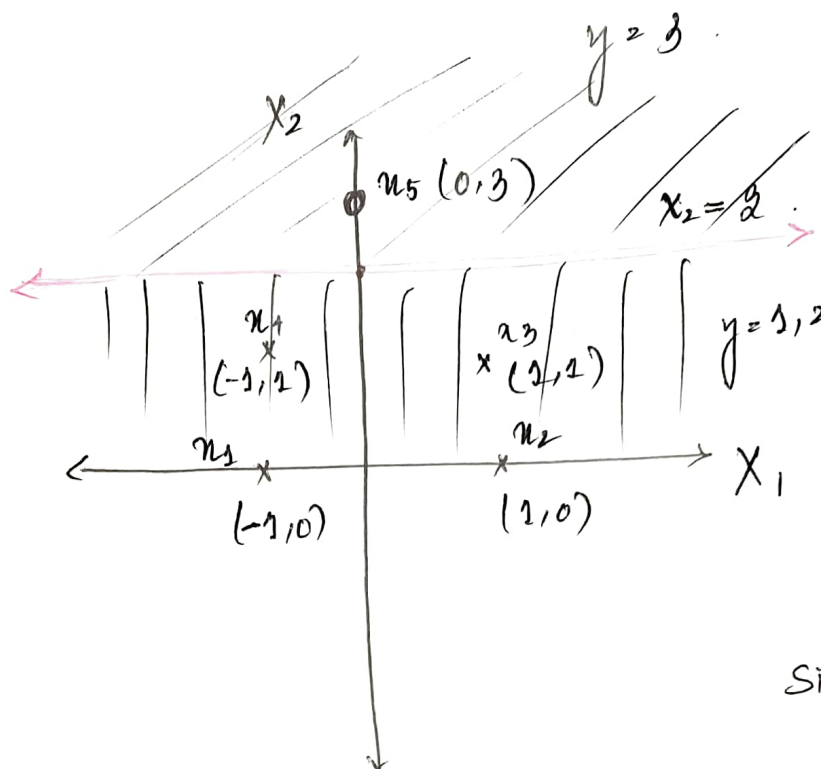
$$\bar{w}^6 = \bar{w}^5 = \begin{bmatrix} 0 \\ -2 \\ 0 \end{bmatrix} = \bar{w}^*$$

$$\begin{aligned}
 h_{\bar{w}}(x) &= \text{sign}(\bar{w}_0 + \bar{w}_1 x_1 + \bar{w}_2 x_2) \\
 &= \text{sign}(-2x_1)
 \end{aligned}$$



5 marks

- (b) Find any linear classification rule that perfectly separates the data points between the two sets: $S_{12} = \{(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4)\}$ and $S_3 = \{(x_5, y_5)\}$. Draw your decision boundary and clearly mark the labels for all the decision regions. You need not use a perceptron algorithm to find the classification rule.



(one of the many possible linear separators)

$$x_2 = 2$$

$$\text{Sign}(2 + 0x_1 - x_2) = \text{Sign}(2 - x_2)$$

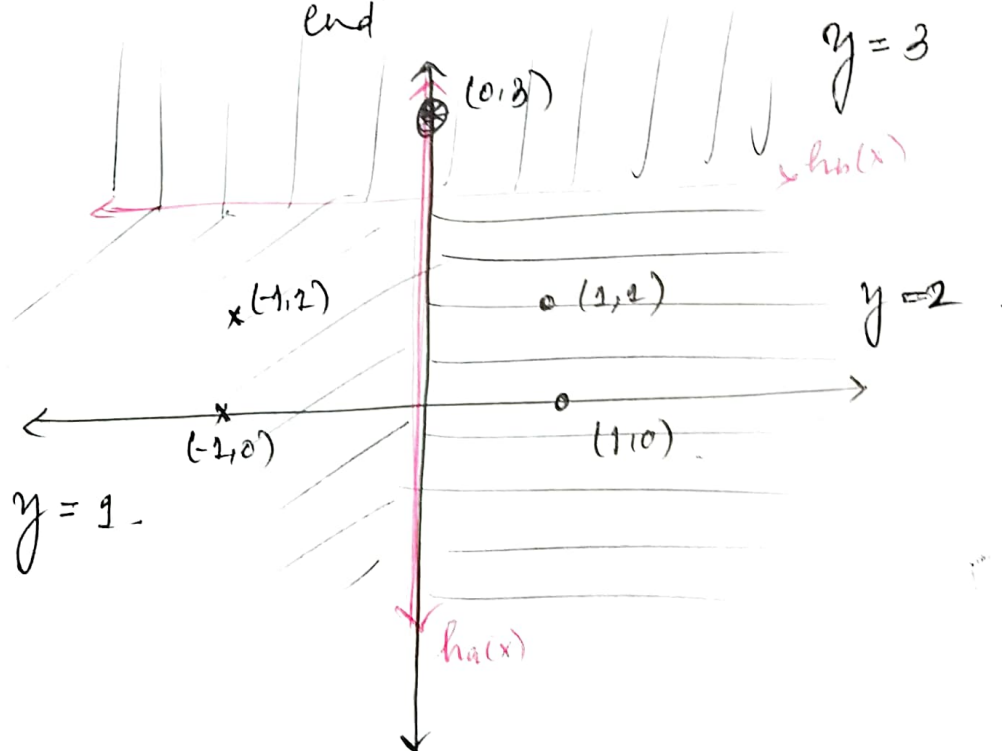
$$\bar{w}^* = \begin{bmatrix} 2 \\ 0 \\ -1 \end{bmatrix} \text{ or } \begin{bmatrix} -2 \\ 0 \\ 1 \end{bmatrix}$$

10 marks

- (c) Explain how to combine parts (a) and (b) to develop a classification rule that given any input $\mathbf{x} \in \mathbb{R}^2$ outputs a label $\hat{y} \in \{1, 2, 3\}$. Your classification rule must achieve perfect classification on the training set. Sketch your decision boundaries in \mathbb{R}^2 and show the labels associated with each decision region.

Ans:

- ① Pass through classifiers in part (b) : $h_b(x)$
- ② If $\hat{y} = -1$ then:
 $y = 3$
 end
- ③ else if $\hat{y} = 1$ then:
 Pass through classifier in part (a) : $h_a(x)$
- ④ If $\hat{y} = -1$:
 $y = 2$
- ⑤ else if $\hat{y} = +1$
 $y = 1$
 end



10 marks

- (d) Suppose we wish to implement a multi-class logistic regression model for classifying the training set \mathcal{D} . Let $\Omega = \{\mathbf{w}(1), \mathbf{w}(2), \mathbf{w}(3)\}$ denote the model parameters of your logistic regression model where $\mathbf{w}(i) \in \mathbb{R}^3$ is the weight vector associated with class label $y = i$. Given an input data vector $\mathbf{x} = (x_1, x_2)^T$ the model outputs is a probability vector:

$$\hat{p}_{\Omega}(i|\mathbf{x}) = \frac{e^{[\mathbf{w}^T(i) \cdot \tilde{\mathbf{x}}]}}{\sum_{j=1}^3 e^{[\mathbf{w}^T(j) \cdot \tilde{\mathbf{x}}]}}, \quad i = 1, 2, 3$$

where $\tilde{\mathbf{x}} = (x_0 = 1, x_1, x_2)^T \in \mathbb{R}^3$ is the augmented vector of \mathbf{x} as discussed in class. We assume a standard log-loss function for the training error, i.e.,

$$E_{\text{in}}(\Omega) = \frac{1}{5} \sum_{n=1}^5 e_n(\Omega), \quad e_n(\Omega) = -\log \hat{p}_{\Omega}(y_n | \mathbf{x}_n)$$

Assuming that we select $\mathbf{w}(1) = (1, 0, 0)^T$, $\mathbf{w}(2) = (0, 1, 0)^T$ and $\mathbf{w}(3) = (0, 0, 1)^T$ numerically evaluate $\nabla_{\mathbf{w}(j)} \{e_1(\Omega)\}$ for $j = 1, 2, 3$.

$$\begin{aligned} \nabla_{\mathbf{w}(j)} \{e_1(\Omega)\} &= \nabla_{\mathbf{w}(j)} (-\log p_{\Omega}(y_n = 1 | \mathbf{x}_n)) \\ &= \nabla_{\mathbf{w}(j)} \left[-\log \left(\frac{e^{\mathbf{w}^T(1) \tilde{\mathbf{x}}_n}}{\sum_{k=1}^3 e^{\mathbf{w}^T(k) \tilde{\mathbf{x}}_n}} \right) \right] \\ &= -\nabla_{\mathbf{w}(j)} (\mathbf{w}^T(1) \tilde{\mathbf{x}}_n) + \frac{\nabla_{\mathbf{w}(j)} \left(\sum_{k=1}^3 e^{\mathbf{w}^T(k) \tilde{\mathbf{x}}_n} \right)}{\sum_{k=1}^3 e^{\mathbf{w}^T(k) \tilde{\mathbf{x}}_n}} \\ &= \cancel{-\tilde{\mathbf{x}}_n \delta(1=j)} + \frac{\tilde{\mathbf{x}}_n e^{\mathbf{w}^T(j) \tilde{\mathbf{x}}_n}}{\sum_{k=1}^3 e^{\mathbf{w}^T(k) \tilde{\mathbf{x}}_n}} \end{aligned}$$

for $n=1; y_1=1$;
for $j=1$;

$$\begin{aligned} \nabla_{\mathbf{w}(1)} \{e_1(\Omega)\} &= \cancel{-\tilde{\mathbf{x}}_1} + \frac{\tilde{\mathbf{x}}_1 e^{\mathbf{w}^T(1) \tilde{\mathbf{x}}_1}}{\sum_{k=1}^3 e^{\mathbf{w}^T(k) \tilde{\mathbf{x}}_1}} \\ &= \frac{-\tilde{\mathbf{x}}_1 + \tilde{\mathbf{x}}_1 e^{\mathbf{w}^T(1) \tilde{\mathbf{x}}_1}}{\sum_{k=1}^3 e^{\mathbf{w}^T(k) \tilde{\mathbf{x}}_1}}, \quad \leftarrow (1) \end{aligned}$$

[continue part (d) here]

for $j=2$:

$$\nabla_{w(2)} \{e_1(\omega)\} = \frac{\tilde{x}_1 e^{\tilde{w}^T(2)\tilde{x}_1}}{\sum_{k=1}^3 e^{\tilde{w}^T(k)\tilde{x}_1}} \quad (2)$$

for $j=3$:

$$\nabla_{w(3)} \{e_1(\omega)\} = \frac{\tilde{x}_1 e^{\tilde{w}^T(3)\tilde{x}_1}}{\sum_{k=1}^3 e^{\tilde{w}^T(k)\tilde{x}_1}} \quad (3)$$

Substituting $\tilde{x}_1 = [1 \ -1 \ 0]^T$ and $\tilde{w}(1) = [1 \ 0 \ 0]^T$ in (2)

$$\begin{aligned} \nabla_{w(1)} \{e_1(\omega)\} &= - \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \left(\frac{1 - e^{[1 \ 0 \ 0]^T [1 \ -1 \ 0]^T}}{e^{[1 \ 0 \ 0]^T [1 \ -1 \ 0]^T} + e^{[0 \ 1 \ 0]^T [1 \ -1 \ 0]^T} + e^{[0 \ 0 \ 1]^T [1 \ -1 \ 0]^T}} \right) \\ &= - \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \left(1 - \frac{e^1}{e^0 + e^1 + e^{-1}} \right) \\ &= \begin{pmatrix} -0.334 \\ 0.334 \\ 0 \end{pmatrix} \end{aligned}$$

Sub \tilde{x}_1 and $\tilde{w}(2)$ in (2)

$$\nabla_{w(2)} \{e_1(\omega)\} = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \left(\frac{e^{[0 \ 1 \ 0]^T [1 \ -1 \ 0]^T}}{e^0 + e^1 + e^{-1}} \right) = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \left(\frac{e^{-1}}{e^0 + e^1 + e^{-1}} \right) = \begin{pmatrix} 0.09 \\ -0.09 \\ 0 \end{pmatrix}$$

Sub \tilde{x}_1 and $\tilde{w}(3)$ in (3):

$$\begin{aligned} \nabla_{w(3)} \{e_1(\omega)\} &= \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \left(\frac{e^{[0 \ 0 \ 1]^T [1 \ -1 \ 0]^T}}{e^0 + e^1 + e^{-1}} \right) = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \left(\frac{e^0}{e^0 + e^1 + e^{-1}} \right) = \\ &= \begin{pmatrix} 0.245 \\ -0.245 \\ 0 \end{pmatrix} \end{aligned}$$

5 marks

- (e) For the problem in part (d), find the output of one-step update of the stochastic gradient descent (SGD) algorithm when the selected training example is $n = 1$, and ϵ is selected as the learning rate.

Update with SGD

$$\begin{aligned} \bar{w}(1) &= \bar{w}(1) - \epsilon (\nabla_{\bar{w}(1)} \{e_1(x_2)\}) \\ &= \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} - \epsilon \begin{pmatrix} -0.334 \\ 0.334 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 + \epsilon 0.334 \\ -\epsilon 0.334 \\ 0 \end{pmatrix} \end{aligned}$$

$$\begin{aligned} \bar{w}(2) &= \bar{w}(2) - \epsilon \{ \nabla_{\bar{w}(2)} \{e_1(x_2)\} \} \\ &= \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} - \epsilon \begin{pmatrix} 0.09 \\ -0.09 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 - \epsilon 0.09 \\ 0 + \epsilon 0.09 \\ 0 \end{pmatrix} \end{aligned}$$

$$\begin{aligned} \bar{w}(3) &= \bar{w}(3) - \epsilon \{ \nabla_{\bar{w}(3)} \{e_1(x_2)\} \} \\ &= \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} - \epsilon \begin{pmatrix} 0.245 \\ -0.245 \\ 0 \end{pmatrix} \\ &= \begin{pmatrix} -\epsilon 0.245 \\ +\epsilon 0.245 \\ 1 \end{pmatrix} \end{aligned}$$

2. (40 MARKS) Consider a linear regression model where the training set is specified by

$$\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\},$$

with $\mathbf{x}_i \in \mathbb{R}^{d+1}$ and $y_i \in \mathbb{R}$. We assume that each data vector is in the augmented dimension i.e., $\mathbf{x}_i = (x_{i,0} = 1, x_{i,1}, \dots, x_{i,d})$. We aim to find a weight vector $\mathbf{w} \in \mathbb{R}^{d+1}$ that aims to minimize a **weighted** squared error loss function

$$E_{\text{in}}^{\Lambda}(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \lambda_i \cdot (\mathbf{w}^T \mathbf{x}_i - y_i)^2,$$

where $\Lambda = (\lambda_1, \dots, \lambda_N)^T$ is a pre-specified vector of **non-negative** constants that determine the importance of each sample.

Suppose that \mathbf{w}^* minimizes the weighted squared error loss i.e.,

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} E_{\text{in}}^{\Lambda}(\mathbf{w}). \quad (1)$$

10 marks

- (a) The optimal solution \mathbf{w}^* can be expressed as a solution to the following expression: $\mathbf{M} \cdot \mathbf{w}^* = \mathbf{U} \cdot \mathbf{y}$ where \mathbf{M} and \mathbf{U} are matrices of dimension $(d+1) \times (d+1)$ and $(d+1) \times N$ respectively and $\mathbf{y} = [y_1, \dots, y_N]^T$ is the observation vector. Provide an expression for \mathbf{M} and \mathbf{U} in terms of the following matrices:

$$\mathcal{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \in \mathbb{R}^{N \times (d+1)} \quad \mathcal{L} = \begin{bmatrix} \sqrt{\lambda_1} & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_2} & \dots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\lambda_N} \end{bmatrix} \in \mathbb{R}^{N \times N}.$$

Please note that \mathcal{L} is a diagonal matrix where the j -th diagonal entry is $\sqrt{\lambda_j}$.

Writing in Matrix and Vector form:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathcal{L}(\mathbf{x}\mathbf{w} - \mathbf{y})\|^2$$

$$0 = \nabla_{\mathbf{w}} ((\mathcal{L}(\mathbf{x}\mathbf{w} - \mathbf{y}))^T (\mathcal{L}(\mathbf{x}\mathbf{w} - \mathbf{y})))$$

$$0 = \nabla_{\mathbf{w}} ((\mathbf{x}\mathbf{w} - \mathbf{y})^T \mathcal{L}^T \mathcal{L} (\mathbf{x}\mathbf{w} - \mathbf{y}))$$

$$0 = \nabla_{\mathbf{w}} ((\mathbf{w}^T \mathbf{x}^T - \mathbf{y}^T) \mathcal{L}^T \mathcal{L} (\mathbf{x}\mathbf{w} - \mathbf{y}))$$

$$0 = \nabla_{\mathbf{w}} ((\mathbf{w}^T \mathbf{x}^T \mathcal{L}^T \mathcal{L} - \mathbf{y}^T \mathcal{L}^T \mathcal{L}) (\mathbf{x}\mathbf{w} - \mathbf{y}))$$

$$0 = \nabla_{\mathbf{w}} (\mathbf{w}^T (\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{x}) \mathbf{w} - \mathbf{w}^T (\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{y}) - \mathbf{y}^T \mathcal{L}^T \mathcal{L} \mathbf{y} - \mathbf{y}^T (\mathcal{L}^T \mathcal{L}) \mathbf{x} \mathbf{w})$$

$$0 = 2 \mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{x} \mathbf{w} - \mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{y} - \mathbf{y}^T \mathcal{L}^T \mathcal{L} \mathbf{x} \mathbf{w}$$

[continue part (a) here]

$$2(x^T L^T L x) \bar{w} = 2(x^T L^T L) y$$

$$(x^T L^T L x) \bar{w} = (x^T L^T L) y$$

$$\therefore M = x^T L^T L x \in \mathbb{R}^{(d+1) \times (d+1)}$$

$$U = x^T L^T L \in \mathbb{R}^{(d+1) \times N}$$

=====

10 marks

- (b) Provide an expression for $\nabla_{\mathbf{w}} (E_{in}^A(\mathbf{w}))$ and use it to provide a (full) gradient descent algorithm for numerically computing the optimal solution \mathbf{w}^* . Assume that a constant learning rate of ϵ is used in the algorithm.

For full grad descent:

$$\begin{aligned}
 \nabla_{\mathbf{w}} (E_{in}^A(\mathbf{w})) &= \nabla_{\mathbf{w}} \left(\frac{1}{N} \sum_{i=1}^N \lambda_i (\mathbf{w}^T \mathbf{x}_i - y_i)^2 \right) \\
 &= \frac{1}{N} \sum_{i=1}^N \lambda_i \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{x}_i - y_i)^2 \\
 &= \frac{1}{N} \sum_{i=1}^N \lambda_i (2(\mathbf{w}^T \mathbf{x}_i - y_i) \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{x}_i - y_i)) \\
 &= \frac{2}{N} \sum_{i=1}^N \lambda_i (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i
 \end{aligned}$$

Update expression:

$$\begin{aligned}
 \mathbf{w}(t+1) &= \mathbf{w}(t) - \epsilon \nabla_{\mathbf{w}} (E_{in}^A(\mathbf{w})) \\
 \mathbf{w}(t+1) &= \mathbf{w}(t) - \frac{2\epsilon}{N} \sum_{i=1}^N \lambda_i (\mathbf{w}^T(t) \mathbf{x}_i - y_i) \mathbf{x}_i
 \end{aligned}$$

5 marks

(c) Using gradient analysis in part (b), provide a stochastic gradient descent (SGD) algorithm for numerically computing the optimal solution \mathbf{w}^* . Assume that a constant learning rate of ϵ is used in the algorithm.

- ① Initialize weights at $t=0$ to $\mathbf{w}(0)$
- ~~② for $t=1, 2, \dots$ epochs do:~~ for $t=1, 2, \dots$ epochs do:
 - ③ for $n \sim \text{uniform } \{1, 2, \dots, N\}$ do:
 - ③ Compute gradient $\mathbf{g}_n(t) = \nabla_{\mathbf{w}} \ell(\mathbf{w})$
 $= 2\lambda_n (\mathbf{w}^T \mathbf{x}_n - y_n) \mathbf{x}_n$
 - ④ Update weights $\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \mathbf{g}_n(t)$
 - ⑤ Return final weights

(can have diff stopping criteria -
eg:
no. of epochs
or
 $E_{in} = 0$
etc)

5 marks

(d) List any two advantages of using the SGD algorithm over the solution in part (a)

- ① Does not need to compute $(\mathbf{x}^T \mathbf{x})^{-1}$ of input matrix, saves memory, computationally efficient.
- ② Not sensitive to outliers.
- ③ Does not easily overfit to training data, generalizing well.

10 marks

(e) Suppose we wish to find \mathbf{w}_β^* minimizes the following:

$$\mathbf{w}_\beta^* = \arg \min_{\mathbf{w} \in \mathbb{R}^{d+1}} \{E_{\text{in}}^\Lambda(\mathbf{w}) + \beta \cdot \|\mathbf{w}\|^2\}, \quad (2)$$

where $E_{\text{in}}^\Lambda(\mathbf{w})$ is the weighted squared error loss as in part (a), $\|\mathbf{w}\|^2$ is the squared Euclidean norm of \mathbf{w} and $\beta > 0$ is the regularization constant. Provide an analytical closed-form expression for \mathbf{w}_β^* in terms of \mathcal{X} , \mathcal{L} , \mathbf{y} , the identity matrix, and the constant β .

$$\mathbf{w}_\beta^* = \arg \min_{\mathbf{w}} \{ \|\mathcal{L}(\mathbf{x}\mathbf{w} - \mathbf{y})\|^2 + \beta \|\mathbf{w}\|^2 \}.$$

$$\begin{aligned} 0 &= \nabla_{\mathbf{w}} \left((\mathcal{L}(\mathbf{x}\mathbf{w} - \mathbf{y}))^T (\mathcal{L}(\mathbf{x}\mathbf{w} - \mathbf{y})) + \beta \mathbf{w}^T \mathbf{w} \right) \\ &= \nabla_{\mathbf{w}} \left(\mathbf{w}^T (\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{x}) \mathbf{w} - \mathbf{w}^T (\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{y}) - \mathbf{y}^T \mathcal{L}^T \mathcal{L} \mathbf{x} \mathbf{w} \right. \\ &\quad \left. + \mathbf{y}^T \mathcal{L}^T \mathcal{L} \mathbf{y} + \beta \mathbf{w}^T \mathbf{w} \right) \end{aligned}$$

$$0 = 2(\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{x}) \mathbf{w}^* - 2(\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{y}) + 2\beta \mathbf{w}^*$$

$$(\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{x}) (\mathbf{w}) + \beta \mathbf{w} = (\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{y})$$

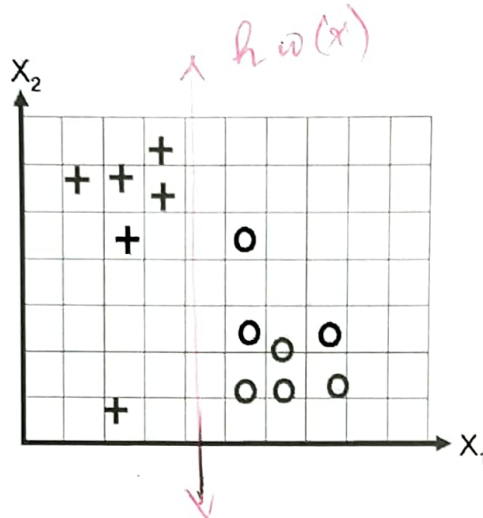
$$(\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{x}) \mathbf{w} + (\beta) (\mathbf{w}^*) = (\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{y})$$

$$(\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{x} + \beta \mathbf{I}) (\mathbf{w}^*) = (\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{y})$$

$$\underline{\underline{\mathbf{w}^* = (\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{x} + \beta \mathbf{I})^{-1} (\mathbf{x}^T \mathcal{L}^T \mathcal{L} \mathbf{y})}}$$

20 marks

3. Consider a binary linear classification problem where $\mathbf{x} \in \mathbb{R}^2$ and $y \in \{-1, +1\}$. We illustrate the training dataset below. The '+' label refers to $y = +1$ and the 'o' label refers to $y = -1$. We would like to construct a classifier $h_{\mathbf{w}}(\mathbf{x}) = \text{sign}(w_0 + w_1x_1 + w_2x_2)$ where $\text{sign}(\cdot)$ is the *sign* function as discussed in class.



In the figure above, the adjacent vertical (and horizontal) lines are 1 unit apart from each other. Assume that the training points are above are $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ (with $N = 13$). We consider the classification loss

$$L(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \mathbb{I}(y_i \neq h_{\mathbf{w}}(\mathbf{x}_i))$$

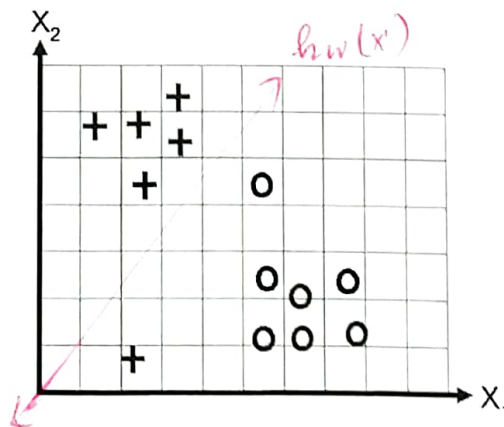
where \mathbb{I} denotes the indicator function.

2 marks

- (a) Draw a decision boundary in the figure above that achieves zero training error.

6 marks

- (b) Suppose that we attempt to minimize the following loss function over \mathbf{w} : $J(\mathbf{w}) = L(\mathbf{w}) + \lambda w_0^2$, where $\lambda = 10^7$ is a huge constant. Sketch a possible decision boundary in the figure below. How many points are mis-classified.



here:

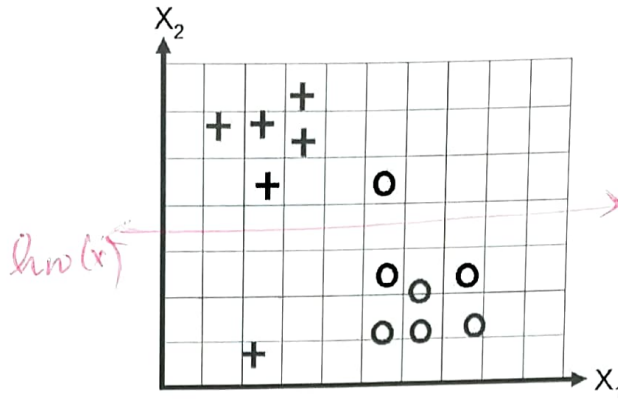
when $\lambda \rightarrow \infty$
 $w_0 \rightarrow 0$

(i.e. decision boundary
 is ~~not~~ passing
 through origin)

1 point is
 misclassified

6 marks

- (c) Suppose that we attempt to minimize the following loss function over \mathbf{w} : $J(\mathbf{w}) = L(\mathbf{w}) + \lambda w_1^2$, where $\lambda = 10^7$ is a huge constant. Sketch a possible decision boundary in the figure below. How many points are mis-classified.

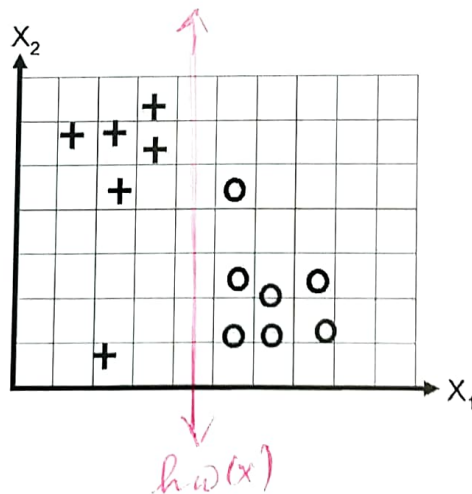


here:
when $\lambda \rightarrow \infty$
 $w_1 \rightarrow 0$
(decision boundary
is parallel to
 x_1)

2 points
misclassified.

6 marks

- (d) Suppose that we attempt to minimize the following loss function over \mathbf{w} : $J(\mathbf{w}) = L(\mathbf{w}) + \lambda w_2^2$, where $\lambda = 10^7$ is a huge constant. Sketch a possible decision boundary in the figure below. How many points are mis-classified.



here:
when $\lambda \rightarrow \infty$
 $w_2 \rightarrow 0$
(decision boundary
parallel to
 x_2)

No points
misclassified.