

FIRST NAME: \_\_\_\_\_ LAST NAME: \_\_\_\_\_

STUDENT NUMBER: \_\_\_\_\_

---

**ECE 1508F — Introduction to Statistical Learning  
Midterm Examination**

Monday November 6<sup>th</sup>, 2017  
1:10 p.m. – 3:00 p.m.

Instructor: Ashish Khisti

*Model Answer & grading scheme*

**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 the 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!

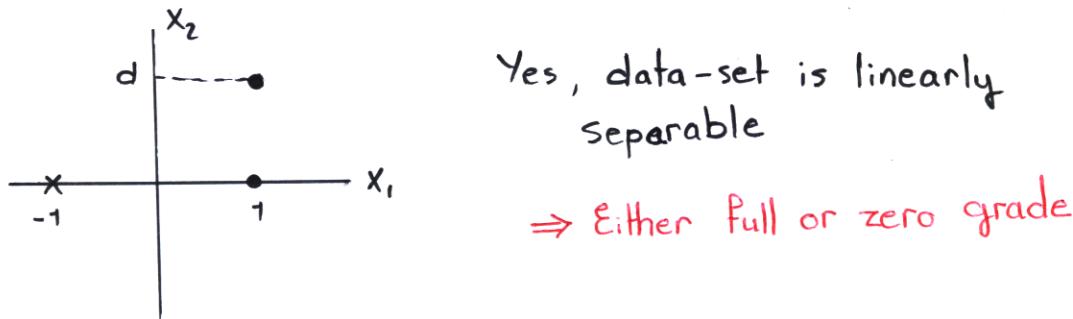
1. (15 MARKS) Consider a binary linear classification where the data points are two dimensional, i.e.,  $\mathbf{x} \in \mathbb{R}^2$  and the labels  $y \in \{-1, 1\}$ . The training set consists of the following points:

$$\begin{aligned}\mathbf{x}_1 &= (1, 0)^T, & y_1 &= +1 \\ \mathbf{x}_2 &= (-1, 0)^T, & y_2 &= -1 \\ \mathbf{x}_3 &= (1, d)^T, & y_3 &= +1\end{aligned}$$

Assume that  $d > 0$  is some constant. For the perceptron algorithm please treat the points that fall exactly on the decision boundary as mistakes and do the update accordingly. Assume that the initial weight is the all-zero vector.

1 mark

- (a) Is the data-set linearly separable? Sketch the points in the  $x_1 - x_2$  plane and show the labels.



6 marks

- (b) Run the perceptron algorithm in the following order  $(1, 2, 3), (1, 2, 3), \dots$  until it converges. Show the output of the perceptron algorithm in each step, sketch the resulting decision boundary, and find the resulting margin (distance of the decision boundary to the nearest training point).

$$\mathbf{w}^{(0)} = [0 \ 0 \ 0]^T \quad \begin{array}{l} \Rightarrow 1 \text{ mark per step} \\ \Rightarrow \text{Not accounting for bias or not initializing to zeros} \end{array} \rightarrow -3$$

$$\text{Step 1 : } \hat{y}_1 = \mathbf{w}^{(0)T} \mathbf{x}_1 = [0 \ 0 \ 0] \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = 0 \rightarrow \text{misclassified}$$

$$\mathbf{w}^{(1)} \leftarrow \mathbf{w}^{(0)} + y_1 \mathbf{x}_1 \implies \mathbf{w}^{(1)} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$$

$$\text{Step 2 : } \hat{y}_2 = \mathbf{w}^{(1)T} \mathbf{x}_2 = [1 \ 1 \ 0] \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} = 0 \rightarrow \text{misclassified}$$

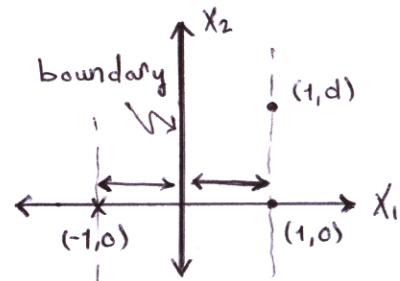
$$\mathbf{w}^{(2)} \leftarrow \mathbf{w}^{(1)} + y_2 \mathbf{x}_2 \implies \mathbf{w}^{(2)} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix}$$

$$\text{Step 3 : } \hat{y}_3 = \mathbf{w}^{(2)T} \mathbf{x}_3 = [0 \ 2 \ 0] \begin{bmatrix} 1 \\ 1 \\ d \end{bmatrix} = 2 \Rightarrow +ve \rightarrow \text{correct}$$

$$\text{Step 4 : } \hat{y}_1 = \mathbf{w}^{(2)T} \mathbf{x}_1 = 2 \Rightarrow +ve \rightarrow \text{correct}$$

$$\text{Step 5 : } \hat{y}_2 = \mathbf{w}^{(2)T} \mathbf{x}_2 = -2 \Rightarrow -ve \rightarrow \text{correct}$$

$$\text{Margin} = 1$$



7 marks

- (c) Repeat part (b) when the order is  $(3, 2, 1), (3, 2, 1, \dots)$ . Also comment on what the decision boundary approaches when  $d \rightarrow \infty$

$$\omega^{(c)\top} = [0 \ 0 \ 0] \quad \begin{array}{l} \Rightarrow 1 \text{ mark per step, 2 marks for margin, 1 for } d \rightarrow \infty \\ \Rightarrow \text{Not accounting for bias or not initializing to zero} \rightarrow -3 \end{array}$$

Step 1:  $\hat{y}_3 = \omega^{(c)\top} x_3 = [0 \ 0 \ 0] \begin{bmatrix} 1 \\ 1 \\ d \end{bmatrix} = 0 \rightarrow \text{misclassified}$

$$\omega^{(1)} \leftarrow \omega^{(c)} + y_3 x_3 = \begin{bmatrix} 1 \\ 1 \\ d \end{bmatrix}$$

Step 2:  $\hat{y}_2 = \omega^{(1)\top} x_2 = [1 \ 1 \ d] \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} = 0 \rightarrow \text{misclassified}$

$$\omega^{(2)} \leftarrow \omega^{(1)} + y_2 x_2 = \begin{bmatrix} 1 \\ 1 \\ d \end{bmatrix} - \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 2 \\ d \end{bmatrix}$$

Step 3:  $\hat{y}_1 = \omega^{(2)\top} x_1 = [0 \ 2 \ d] \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = 2 \Rightarrow +ve \rightarrow \text{correct}$

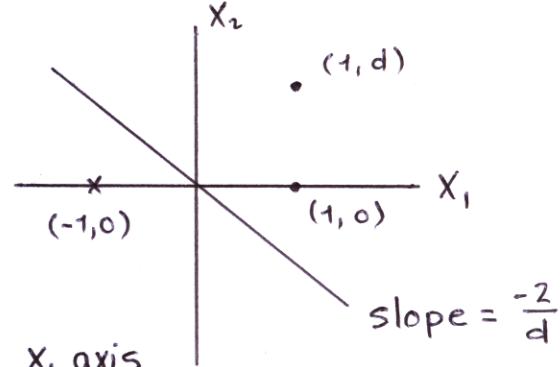
Step 4:  $\hat{y}_3 = \omega^{(2)\top} x_3 = [0 \ 2 \ d] \begin{bmatrix} 1 \\ 1 \\ d \end{bmatrix} = +2 \Rightarrow +ve \rightarrow \text{correct}$

∴ Boundary is  $2x_1 + dx_2 = 0$

$$x_2 = \frac{-2}{d} x_1$$

$$\text{Margin} = \frac{|(2 \times 1) + (d \times 0)|}{\sqrt{(2)^2 + (d)^2}} = \frac{2}{\sqrt{4 + d^2}}$$

As  $d \rightarrow \infty$ : slope  $\rightarrow 0$  ∴ Boundary is  $x_1$  axis ( $x_2 = 0$ )



1 mark

- (d) Between parts (b) and (c) which solution will be preferred? Briefly explain.

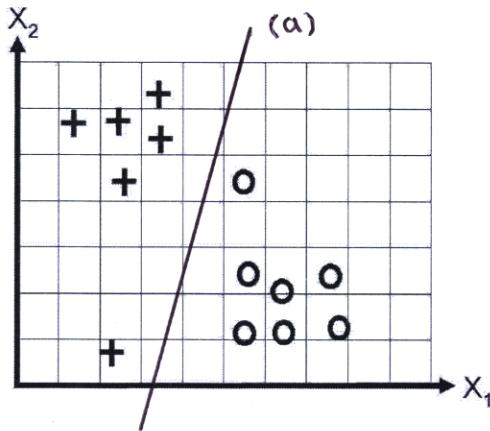
(b) is preferred as it has bigger margin

$$\left( \frac{2}{\sqrt{4+d^2}} < 1 \right)$$

⇒ Either full or zero grade

10 marks

2. 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

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

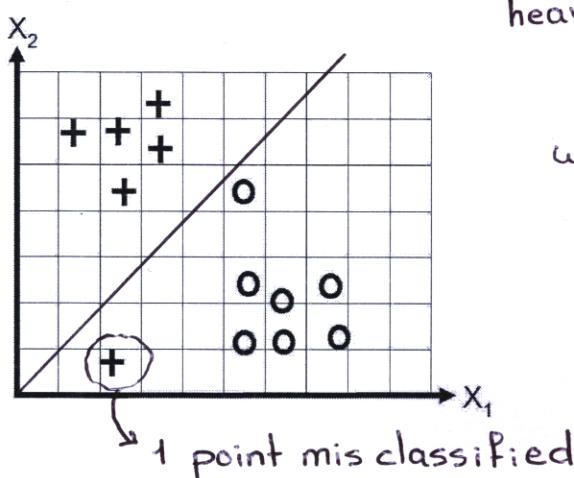
where  $\mathbb{I}$  denotes the indicator function.

2 marks

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

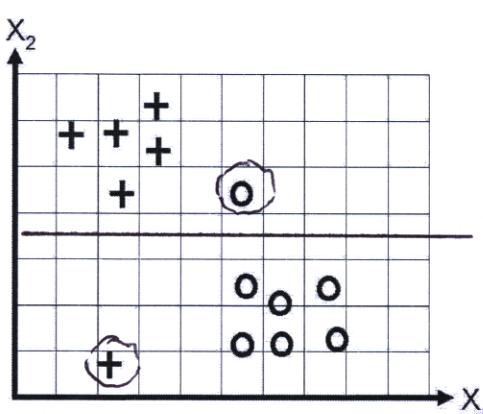
4 marks

- (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.



2 marks

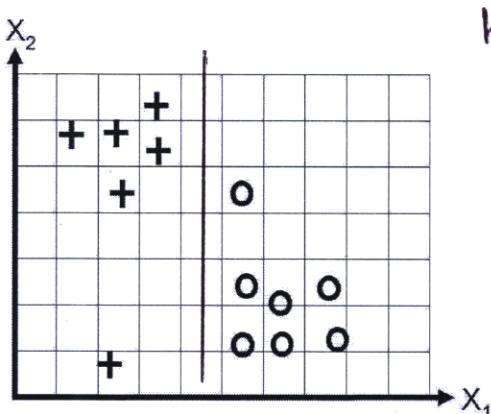
- (c) Suppose that we attempt to minimize the following loss function over  $\mathbf{w}$  :  $J(\mathbf{w}) = L(\mathbf{w}) + \lambda w_1^2$ , where  $\lambda = 10^7$  is a huge constant. Sketch a possible decision boundary in the figure below. How many points are mis-classified.



heavily penalizing  $w_1$   
 $\therefore w_1 \rightarrow 0$   
 $w_0 + w_2 x_2 = 0$  is  
 the boundary  
 2 points misclassified

2 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.



heavily penalizing  $w_2$   
 $\therefore w_2 \rightarrow 0$   
 $w_0 + w_1 x_1 = 0$  is  
 the boundary  
 No points misclassified

⇒ For parts (b), (c) & (d) , half the question grade for sketching & the other half for number of misclassified points

- 25 marks 3. Suppose that a data vector  $\mathbf{y} \in \mathbb{R}^n$  and a matrix  $\mathbf{A} \in \mathbb{R}^{n \times p}$  are given. Assume that  $n > p$  and that  $\mathbf{A}^T \mathbf{A}$  is an  $p \times p$  invertible matrix. We would like to approximate  $\mathbf{y}$  using a vector of the form  $\hat{\mathbf{y}}_{\mathbf{w}} = \mathbf{A}\mathbf{w}$ .

- 2 marks (a) Let  $\mathbf{w}_{LS} = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{y}_{\mathbf{w}}\|^2$ , be the least square estimate of  $\mathbf{y}$ . Provide an expression for  $\mathbf{w}_{LS}$  and the associated  $\hat{\mathbf{y}}_{LS} = \mathbf{A}\mathbf{w}_{LS}$ . No calculation is required in this part.

$$\mathbf{w}_{LS} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} \quad \Rightarrow 1 \text{ point for each}$$

$$\hat{\mathbf{y}}_{LS} = \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$

- 4 marks (b) Show that for any  $\mathbf{w} \in \mathbb{R}^p$  the following identity holds:

$$\|\mathbf{y} - \hat{\mathbf{y}}_{\mathbf{w}}\|^2 = \|\mathbf{y} - \hat{\mathbf{y}}_{LS}\|^2 + \|\hat{\mathbf{y}}_{\mathbf{w}} - \hat{\mathbf{y}}_{LS}\|^2$$

No calculation is needed for this part, just a clear geometric argument and an accompanying figure that uses properties of  $\hat{\mathbf{y}}_{LS}$  is sufficient.

$\hat{\mathbf{y}}_{LS}$  is the projection of  $\mathbf{y}$  onto  $\text{col-span}\{\mathbf{A}\}$

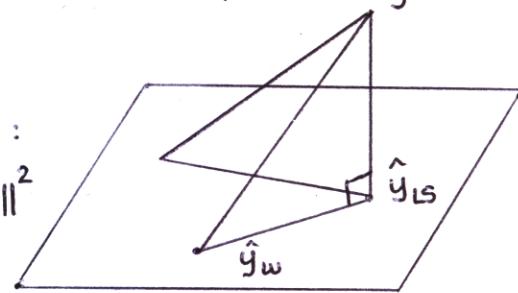
$\therefore \mathbf{y} - \hat{\mathbf{y}}_{LS}$  is  $\perp$  to  $\hat{\mathbf{y}}_{\mathbf{w}} - \hat{\mathbf{y}}_{LS}$

Using pythagoras theorem:

$$\|\mathbf{y} - \hat{\mathbf{y}}_{\mathbf{w}}\|^2 = \|\mathbf{y} - \hat{\mathbf{y}}_{LS}\|^2 + \|\hat{\mathbf{y}}_{\mathbf{w}} - \hat{\mathbf{y}}_{LS}\|^2$$

$\Rightarrow 2$  points for sketch

$\Rightarrow 2$  points for explanation



- 6 marks (c) For parts (c) and (d) of this question assume that  $\mathbf{A}^T \mathbf{A}$  is a diagonal matrix and  $\lambda > 0$  is a constant. Let  $a_1, a_2, \dots, a_p$  be the elements on the diagonal of  $\mathbf{A}^T \mathbf{A}$ , assumed to be non-zero. Let  $\mathbf{w}^* = \arg \min_{\mathbf{w}} \{\|\mathbf{y} - \hat{\mathbf{y}}_{\mathbf{w}}\|^2 + \lambda \|\mathbf{w}\|^2\}$ . Using part (b) show that

$$w_i^* = \arg \min_w \{a_i(w - w_{LS,i})^2 + \lambda w^2\}, \quad i = 1, 2, \dots, p$$

where  $w_i^*$  and  $w_{LS,i}$  denote the  $i$ -th element of  $\mathbf{w}^*$  and  $\mathbf{w}_{LS}$  respectively. Hence express  $w_i^*$  in terms of  $w_{LS,i}$ ,  $a_i$  and  $\lambda$ .

$\Rightarrow 3$  points for  $\mathbf{w}^*$

$$\left\{ \begin{array}{l} \mathbf{w}^* = \arg \min_{\mathbf{w}} \{\|\mathbf{y} - \hat{\mathbf{y}}_{\mathbf{w}}\|^2 + \lambda \|\mathbf{w}\|^2\} \text{ (using part (b))} \\ = \arg \min_{\mathbf{w}} \{\|\mathbf{y} - \hat{\mathbf{y}}_{LS}\|^2 + \|\hat{\mathbf{y}}_{\mathbf{w}} - \hat{\mathbf{y}}_{LS}\|^2 + \lambda \|\mathbf{w}\|^2\} \text{ (\|y - y_Ls\|^2 is independant of w)} \\ = \arg \min_{\mathbf{w}} \{\|\hat{\mathbf{y}}_{\mathbf{w}} - \hat{\mathbf{y}}_{LS}\|^2 + \lambda \|\mathbf{w}\|^2\} \\ = \arg \min_{\mathbf{w}} \{(\mathbf{w} - \mathbf{w}_{LS})^T \mathbf{A}^T \mathbf{A} (\mathbf{w} - \mathbf{w}_{LS}) + \lambda \mathbf{w}^T \mathbf{w}\} \end{array} \right.$$

$\Rightarrow 3$  points for getting  $w_i^*$

$$\left\{ \begin{array}{l} \because \mathbf{A}^T \mathbf{A} \text{ is a diagonal matrix} \\ \therefore w_i^* = \arg \min_{\mathbf{w}} \{a_i(w_i - w_{LS,i})^2 + \lambda w_i^2\} = \arg \min_{\mathbf{w}} \{h(w_i)\} \\ h'(w_i) = 0 \rightarrow 2a_i(w_i - w_{LS,i}) + 2\lambda w_i = 0 \\ \therefore w_i = \frac{a_i}{a_i + \lambda} w_{LS,i} \end{array} \right.$$

4 marks

(d) Suppose we wish to compute:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left\{ \|\mathbf{y} - \hat{\mathbf{y}}_{\mathbf{w}}\|^2 + \lambda \sum_{i=1}^p |w_i| \right\}$$

where  $|\cdot|$  is the absolute value. Using the method analogous to part (c) express  $w_i^*$  (the  $i$ -th component of  $\mathbf{w}^*$ ) in terms of  $w_{LS,i}$ ,  $a_i$  and  $\lambda$ . Hint: Consider the function  $f(x) = (x - c)^2 + \eta|x|$ , where  $\eta > 0$ . The minimizing value of  $f(x)$ , say  $x_0$ , can be expressed as follows: If  $c > \eta/2$  then  $x_0 = c - \eta/2$ , if  $c < -\eta/2$  then  $x_0 = c + \eta/2$  and if  $|c| \leq \eta/2$  then  $x_0 = 0$ .

$$\begin{aligned} \text{from part (c)} : w_i^* &= \operatorname{argmin} \{ a_i (w_i - w_{LS,i})^2 + \lambda |w_i| \} \\ &= \operatorname{argmin} \{ (w_i - w_{LS,i})^2 + \frac{\lambda}{a_i} |w_i| \} \end{aligned}$$

(Using the hint given)

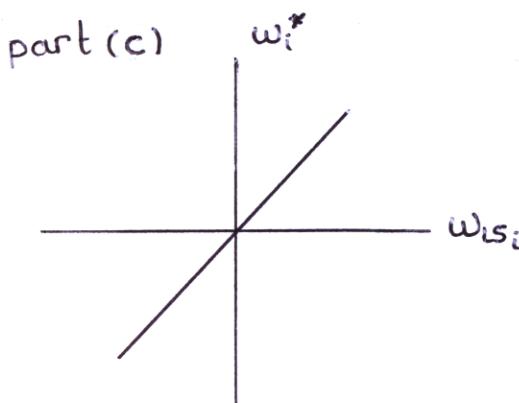
$$= \begin{cases} w_{LS,i} - \frac{\lambda}{2a_i}, & w_{LS,i} > \frac{\lambda}{2a_i} \\ w_{LS,i} + \frac{\lambda}{2a_i}, & w_{LS,i} < -\frac{\lambda}{2a_i} \\ 0, & \text{otherwise} \end{cases}$$

⇒ 2 points for putting  $w_i^*$  in the correct format

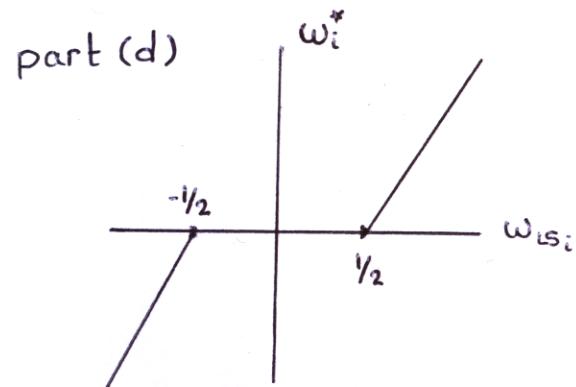
⇒ 2 points for using the hint

4 marks

(e) Sketch  $w_i^*$  as a function of  $w_{LS,i}$  for  $\lambda = 1$  and  $a_i = 1$  for both parts (c) and (d). Also explain qualitatively how does your answer in the two cases differ for large values of  $\lambda$ ?



for large  $\lambda$ :  $w_i^* = \frac{a_i}{a_i + \lambda} w_{LS,i}$   
 $\lambda \rightarrow \infty \quad w_i^* \rightarrow 0$



for large  $\lambda$ :  $w_i^* = 0$

⇒ 1 point for sketch & 1 point for large  $\lambda$  in each part

5 marks

- (f) Suppose that we have  $N$  training examples:  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  where  $\mathbf{x}_i \in \mathbb{R}^p$  and we wish to minimize the following loss function:

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i^T \mathbf{w} - y_i)^2 + \frac{\lambda}{N} \sum_{j=1}^p |w_j|,$$

where  $\mathbf{w} = (w_1, \dots, w_p)^T$ . Provide a SGD update rule for minimizing  $J(\mathbf{w})$ . In each step assume that one example is selected at random in each update, and the gradient from the regularization term is added in each step.

$$\underline{\mathbf{w}}(t+1) \leftarrow \underline{\mathbf{w}}(t) - \eta_t \nabla_{\underline{\mathbf{w}}(t)} e_n$$

$$e_n = (\mathbf{x}_n^T \underline{\mathbf{w}} - y_n)^2 + \lambda \sum_{j=1}^p |w_j|$$

$$\frac{\partial e_n}{\partial w_j} = 2(\mathbf{x}_n^T \underline{\mathbf{w}} - y_n) x_j + \lambda \text{sign}(w_j)$$

$$\nabla_{\underline{\mathbf{w}}(t)} e_n = 2(\mathbf{x}_n^T \underline{\mathbf{w}}(t) - y_n) \underline{\mathbf{x}}_n + \lambda \text{sign}(\underline{\mathbf{w}}(t))$$

$\Rightarrow$  2 points for update rule (knowing that  $N=1$ )  
 $\Rightarrow$  3 points for derivatives

10 marks

4. Consider a logistic regression binary classification model that given  $\mathbf{x} \in \mathbb{R}^2$  outputs:

$$P_{\mathbf{w}}(y=1|\mathbf{x}) = \frac{1}{1+e^{-(w_1x_1+w_2x_2)}}, \quad P_{\mathbf{w}}(y=-1|\mathbf{x}) = \frac{1}{1+e^{(w_1x_1+w_2x_2)}}.$$

Assume that the bias term in the model is set to zero for simplicity in this problem. Suppose that we have only two training points:

$$\begin{aligned}\mathbf{x}_1 &= (1, 1), \quad y_1 = 1 \\ \mathbf{x}_2 &= (1, 0), \quad y_2 = -1\end{aligned}$$

We intend to select  $\mathbf{w}$  that minimizes the following regularized loss function:

$$J(\mathbf{w}) = -\sum_{i=1}^2 \log P_{\mathbf{w}}(y_i|\mathbf{x}_i) + \lambda \|\mathbf{w}\|^2$$

5 marks

- (a) Suppose that  $\lambda = 0$ . What will the optimal choice of  $\mathbf{w}$  be? What will the resulting value of  $J(\mathbf{w})$  be?

$$\begin{aligned}\because w_0 = 0, \lambda = 0 : J(w) &= -\log\left(\frac{1}{1+e^{-(w_1+w_2)}}\right) - \log\left(\frac{1}{1+e^{w_1}}\right) \\ &= \log(1+e^{-(w_1+w_2)}) + \log(1+e^{w_1})\end{aligned}$$

To minimize  $J(w) \Rightarrow w_1 = -K, w_2 = 2K \quad \& \quad K \rightarrow \infty$   
 $\therefore J(w) = 0$

$\Rightarrow$  2 points for writing eq. of  $J(w)$

$\Rightarrow$  2 points for  $w \rightarrow \infty$ , 1 point for  $J(w) = 0$

5 marks

- (b) Suppose that  $\lambda$  is a large constant such that it only suffices to consider  $\|\mathbf{w}\| \ll 1$  when minimizing  $J(\mathbf{w})$ . In this case we can approximate

$$\log(1+e^{-y_i \cdot \mathbf{w}^T \mathbf{x}_i}) \approx \log 2 - \frac{1}{2} y_i \cdot \mathbf{w}^T \mathbf{x}_i$$

Assuming that the above approximation is exact find  $\mathbf{w}$  that minimizes  $J(\mathbf{w})$ .

$$\begin{aligned}J(w) &= \log(1+e^{-(w_1+w_2)}) + \log(1+e^{w_1}) + \lambda(w_1^2 + w_2^2) \\ &= \log 2 - \frac{1}{2}(w_1+w_2) + \log 2 + \frac{1}{2}w_1 + \lambda(w_1^2 + w_2^2)\end{aligned}$$

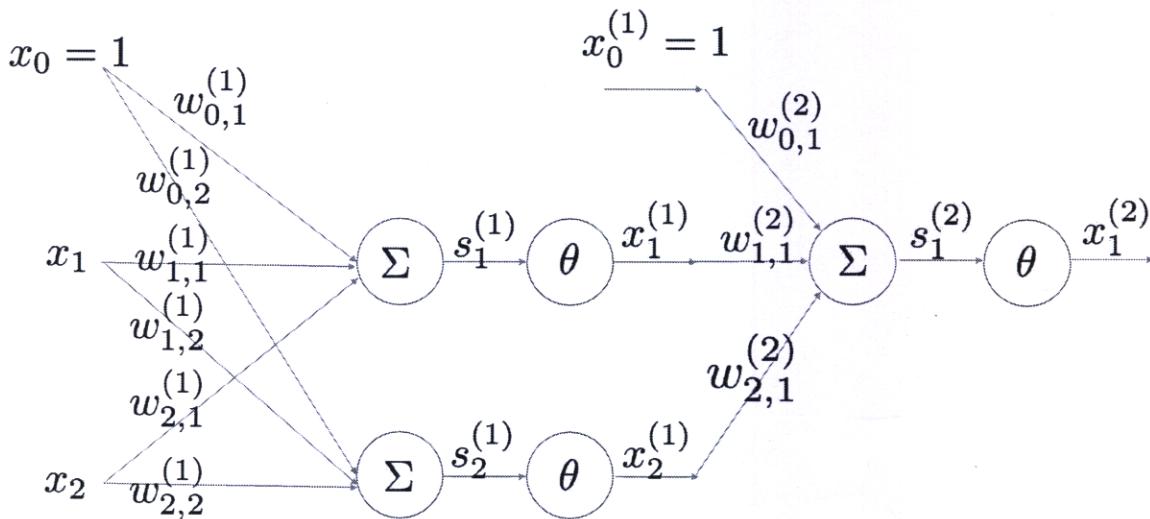
$$\begin{aligned}\frac{\partial J(w)}{\partial w_1} &= -\frac{1}{2} + \frac{1}{2} + 2\lambda w_1 = 0 \quad \rightarrow \quad w_1 = 0 \\ \frac{\partial J(w)}{\partial w_2} &= -\frac{1}{2} + 2\lambda w_2 = 0 \quad \rightarrow \quad w_2 = \frac{1}{4\lambda}\end{aligned}\quad \left. \begin{array}{l} \\ \end{array} \right\} \quad \begin{array}{l} w^* = \begin{bmatrix} 0 \\ \frac{1}{4\lambda} \end{bmatrix} \end{array}$$

$\Rightarrow$  3 points for writing eq. of  $J(w)$

$\Rightarrow$  2 points for getting  $w_1$  &  $w_2$

20 marks

5. Consider a neural network with 2 layers as shown below. Note that  $w_{i,j}^{(l)}$  denotes the weight on the edge between node  $i$  in layer  $l-1$  and node  $j$  in layer  $l$ . The input symbols are denoted by  $x_1$  and  $x_2$  and the output symbol is denoted by  $x_1^{(2)}$ . The symbol  $\Sigma$  denotes summation while  $\theta$  denotes the activation function. Thus  $s_1^{(1)} = w_{0,1}^{(1)}x_0 + w_{1,1}^{(1)}x_1 + w_{2,1}^{(1)}x_2$  and  $x_1^{(1)} = \theta(s_1^{(1)})$ , for example.



5 marks

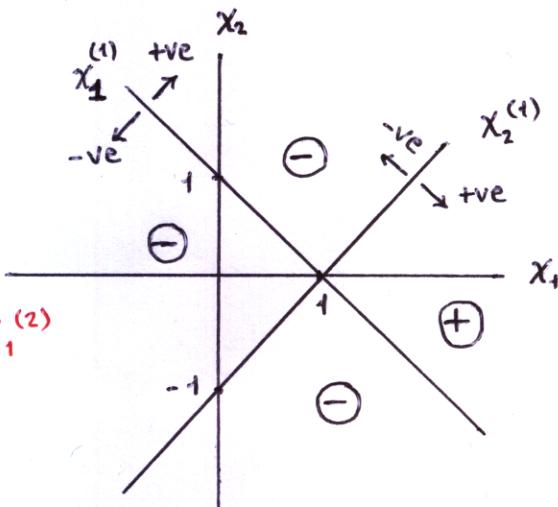
- (a) Suppose that we let  $\theta(s) = \text{sign}(s) \in \{-1, 1\}$ , and we select the following weights:

$$\begin{aligned} w_{0,1}^{(1)} &= -1, & w_{1,1}^{(1)} &= 1, & w_{2,1}^{(1)} &= 1 \\ w_{0,2}^{(1)} &= -1 & w_{1,2}^{(1)} &= 1 & w_{2,2}^{(1)} &= -1 \\ w_{0,1}^{(2)} &= -1.5 & w_{1,1}^{(2)} &= 1 & w_{2,1}^{(2)} &= 1 \end{aligned}$$

In the  $x_1 - x_2$  plane, sketch the regions where  $x_1^{(2)} = +1$  and the regions where  $x_1^{(2)} = -1$

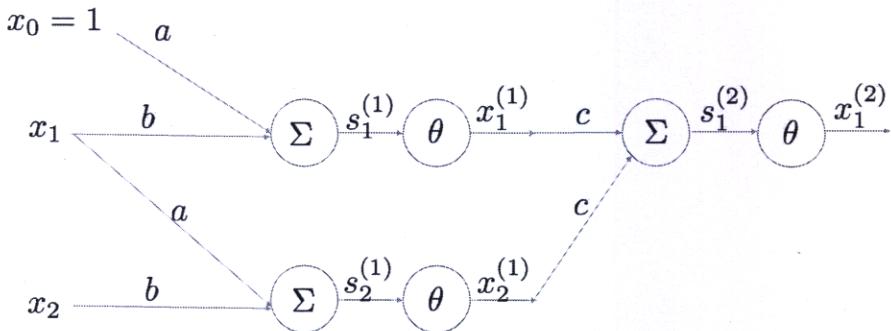
$$\begin{aligned} x_1^{(1)} &= \text{sign}(x_1 + x_2 - 1) \\ x_2^{(1)} &= \text{sign}(x_1 - x_2 - 1) \\ x_1^{(2)} &= \text{AND}(x_1^{(1)}, x_2^{(1)}) \end{aligned}$$

⇒ 3 points for  $x_1^{(1)}$ ,  $x_2^{(1)}$  &  $x_1^{(2)}$   
 ⇒ 2 points for sketch



10 marks

- (b) Suppose that we impose the following restrictions on the weights in the neural network in the previous page:  $w_{0,1}^{(0)} = w_{2,1}^{(1)} = w_{0,1}^{(2)} = 0$  and furthermore  $w_{0,1}^{(1)} = w_{1,2}^{(2)}$ ,  $w_{1,1}^{(1)} = w_{2,2}^{(1)}$ ,  $w_{1,1}^{(2)} = w_{2,1}^{(2)}$ . The resulting neural network, shown below, uses simplified variables  $a = w_{0,1}^{(1)}$ ,  $b = w_{1,1}^{(1)}$  and  $c = w_{1,1}^{(2)}$ . The parameters  $a$ ,  $b$  and  $c$  need to be learned in the network.



Given a training example  $\mathbf{x} = (x_1, x_2)$  with label  $y$ , let the loss function be defined as:

$$e = (x_1^{(2)} - y)^2.$$

Assuming that  $\theta(s)$  be an arbitrary activation function with derivative  $\theta'(s)$  provide expressions for  $\frac{\partial e}{\partial a}$ ,  $\frac{\partial e}{\partial b}$ ,  $\frac{\partial e}{\partial c}$ ,  $\frac{\partial e}{\partial x_1}$  and  $\frac{\partial e}{\partial x_2}$  by using a back-propagation type algorithm.

$$\frac{\partial e}{\partial s_1^{(2)}} = 2(x_1^{(2)} - y)\theta'(s_1^{(2)}) = \delta_1^{(2)}$$

$$\frac{\partial e}{\partial c} = \delta_1^{(2)}(x_1^{(1)} + x_2^{(1)})$$

$$\frac{\partial e}{\partial s_1^{(1)}} = \delta_1^{(2)} \cdot \theta'(s_1^{(1)}) = \delta_1^{(1)}$$

$$\frac{\partial e}{\partial s_2^{(1)}} = \delta_1^{(2)} \cdot \theta'(s_2^{(1)}) = \delta_2^{(1)}$$

$$\frac{\partial e}{\partial a} = \delta_1^{(1)} x_0 + \delta_2^{(1)} x_1$$

$$\frac{\partial e}{\partial x_1} = \delta_1^{(1)} b + \delta_2^{(1)} a$$

$$\frac{\partial e}{\partial b} = \delta_1^{(1)} x_1 + \delta_2^{(1)} x_2$$

$$\frac{\partial e}{\partial x_2} = \delta_2^{(1)} b$$

$\Rightarrow$  1 point for  $\frac{\partial e}{\partial c}$  & 2 points for every other requirement

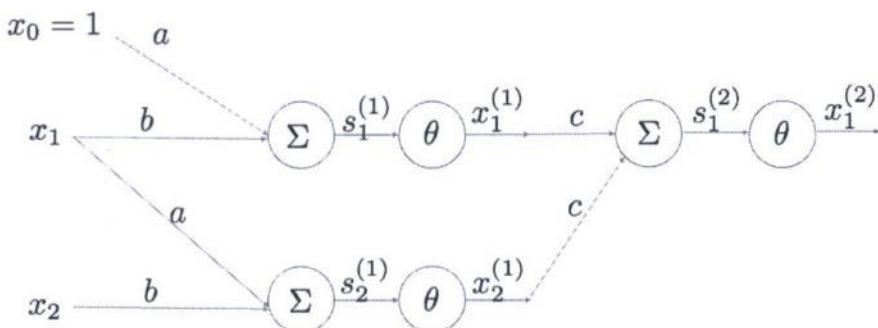
[use this page for extra work for part (b).]



total/10

5 marks

- (c) For the neural network in part (b) (reproduced below) numerically evaluate  $\frac{\partial e}{\partial x_1}$  and  $\frac{\partial e}{\partial x_2}$  in the following case:  $\theta(s) = \max(0, s)$ ,  $x_1 = 2$ ,  $x_2 = 1$ ,  $y = 0$ ,  $a = 1$ ,  $b = -1$ ,  $c = 1$ .



Forward pass :  $x_1^{(2)} = 1$

Backward pass :  $\delta_1^{(2)} = 2 \rightarrow \frac{\partial e}{\partial x_1^{(1)}} = 2$ ,  $\frac{\partial e}{\partial x_2^{(1)}} = 2$   
 (using eq.s  
 from b)

$$\delta_1^{(1)} = 0$$

$$\delta_2^{(1)} = 2$$

$$\therefore \frac{\partial e}{\partial x_1} = a \delta_2^{(1)} + b \delta_1^{(1)} = 2$$

$$\therefore \frac{\partial e}{\partial x_2} = b \delta_2^{(1)} = -2$$

⇒ 1 point for Forward pass  
 2 points for  $\delta_1^{(1)}$  and  $\delta_2^{(1)}$   
 2 points for final answers