

# Problem Set 1

## CS 6375

WarmUp:

Subgradient & More

(a)  $f(x) = \max \{x^2, (x+1)^2\}$  at  $x=1$  and  $x=0.5$

At  $x=1$   $\left\{ \begin{array}{l} x^2 = 1 \\ (x+1)^2 = 4 \end{array} \right.$

$$\therefore f'(x) = 2(x+1)$$

$\therefore$  the subgradient at  $x=1$  is  $2(x+1) = 4$

At  $x=0.5$   $\left\{ \begin{array}{l} (0.5)^2 = 0.25 \\ (0.5+1)^2 = 2.25 \end{array} \right.$

$$\therefore f'(x) = 2(x+1)$$

$\therefore$  the subgradient at  $x=0.5$   $2(x+1) = 3$

(b)  $g(x) = \max \{x^2, \exp(x), 10x\}$  at  $x=-1$  and  $x=0.01$

At  $x=-1$   $\left\{ \begin{array}{l} x^2 = 1 \\ e^x = e^{-1} \\ 10x = -10 \end{array} \right.$

$$f'(x) = 2x$$

$\therefore$  the subgradient of  $f(x)$  at  $x = -1$  is

$$2x = -2$$

At  $x = 0.01$

$$\begin{cases} 2^2 & 100^{-2} \\ e^2 & e^{0.01} \\ 10x & 0.10 \end{cases}$$

$\therefore$  the subgradient, at  $x=0.01$ , is  $f'(x) = 10$

### Problem #1

#### 1. Separability & Feature Vectors

Perception loss using subgradient descent

$$\omega^{(t+1)} = \omega^{(t)} + y_t \cdot \sum_{i=y_t f(x^{(i)}) > 0} y_i x^{(i)}$$

$$b^{(t+1)} = b^{(t)} + y_t \cdot \sum_{i=-y_t f(x^{(i)}) > 0} y_i$$

The  $\omega^{t+1}$  and  $b^{t+1}$  will be different from  $\omega^t$  and  $b^t$  only when the example is misclassified.

If the training set is not linearly separable

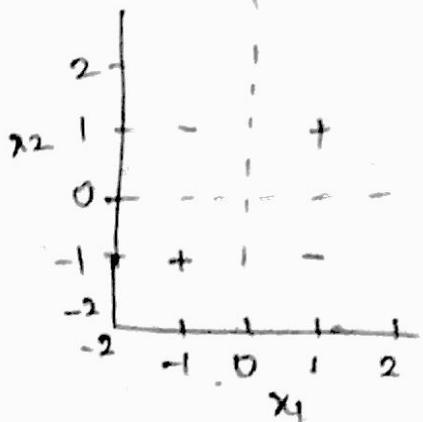
i.e. no hyperplane can divide it into two separate classes. we will always have misclassified points/samples that will try to change  $w^T$  and  $b$  accordingly

Thus we can never find a solution that can separate the training data.

Therefore, the perceptron loss ~~cost~~ may not converge if the training data is not linearly separable.

2.

Data set given



As it is not linearly separable, adding feature vectors will make the data linearly separable.

Using intuition we can see that in the first quarter and third quarter it is +ve and other two quarters it is negative.

So a hyperbole i.e a rectangular hyperbole  
~~will~~ will transform the data to make it  
 linearly separable.

~~rectangular~~  
 right/~~rectangle~~

$$\therefore \underline{xy = c} \quad \text{equation of hyperbole}$$

of the given ~~or~~ feature vectors in the question.

$$(e) \phi(x_1, x_2) = \begin{bmatrix} x_1 x_2 \\ 1 \end{bmatrix}$$

looks the feature vector.

~~is~~ therefore, now we will find  $w^T$  and  $b$ .  
 to check if it is converging or not.

$$w^0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad b^0 = 0$$

$$\text{Taking tuple } t: \begin{matrix} x_1 & x_2 & y \end{matrix} \quad \therefore \phi(x_1, x_2) = \begin{bmatrix} \cdot \\ 1 \end{bmatrix}$$

$$w^0 \phi^t + b^0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \begin{bmatrix} \cdot \\ 1 \end{bmatrix} + 0 = 0.$$

$\therefore$   $y_t = 1 \rightarrow y_t = 1$  ] for all cases

$$w' = w^0 + y_t y_t \phi(x_1, x_2) \quad \left| \begin{array}{l} b' = b^0 + y \\ = 1 \end{array} \right.$$

$$= \begin{bmatrix} 0 \\ 0 \end{bmatrix} + (1) \begin{bmatrix} \cdot \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} \cdot \\ 1 \end{bmatrix}$$

Next data input  $t = \begin{matrix} x_1 & x_2 & y \\ -1 & 1 & -1 \end{matrix}$   $\phi(x_1, x_2) = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$

$$\omega^1 \phi(x_1, x_2) + b^1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} + 1 = 0$$

As it is misclassified.

$$\omega^2 = \omega^1 + y_t y \phi(x_1, x_2)$$

$$= \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 1 \times -1 \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$b^2 = b^1 + y_t y$$

$$= 0 + -1$$

$$b^2 = 0$$

$$\omega^2 = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

Next data input  $t = \begin{matrix} x_1 & x_2 & y \\ 1 & -1 & -1 \end{matrix}$   $\phi(x_1, x_2) = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$

$$\omega^2 \phi(x_1, x_2) + b^2 = \begin{bmatrix} 2 \\ 0 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} + 0 = -\underline{\text{ve}}$$

As it is correctly classified, there is no need to modify weights

Next data input  $t = \begin{matrix} x_1 & x_2 & y \\ 1 & 1 & 1 \end{matrix}$   $\phi(x_1, x_2) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$$\omega^2 \phi(x_1, x_2) + b^2 = \begin{bmatrix} 2 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 0 = +\text{ve}$$

Again it is correctly classified.

Again we start from top.

Next point  $t = -1 \quad x_1 \quad x_2 \quad y$ ,  $\phi(x_1, x_2) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$$w^T \phi(x_1, x_2) + b^2 = \begin{bmatrix} 2 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 0 = \text{+ve}$$

Correctly classified.

Finally the last point to check.

$t = -1 \quad x_1 \quad x_2 \quad y$ ,  $\phi(x_1, x_2) = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$

$$w^T \phi(x_1, x_2) + b^2 = \begin{bmatrix} 2 \\ 0 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} + 0 = \text{-ve}$$

Correctly classified.

∴ The linear separator using ~~feature~~ feature variable is

$$w^T \phi(x_1, x_2) + b = 0$$

where  $w^T = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$ ,  $\phi(x_1, x_2) = \begin{bmatrix} x_1 x_2 \\ 1 \end{bmatrix}$ ,  $b = 0$

## Problem #2 : Perceptron Learning

$$\sum_{i=1}^n \max \{0, -y_i (\omega^T x^{(i)} + b)\}$$

### ① Standard Gradient Descent

In standard gradient descent we minimize the function by moving in the direction where the loss is less.

In other words, we modify the function such that the total error/loss is decreased.

Therefore in Perceptron loss algorithm

To minimize the perceptron loss we use the equations

$$\omega^{(t+1)} = \omega^{(t)} + y_t \cdot \sum_{i=y_i; f(x^{(i)}) > 0} y_i x^{(i)}$$

$$b^{(t+1)} = b^{(t)} + y_t \cdot \sum_{i=y_i; f(x^{(i)}) > 0} y_i$$

We start with  $w^0 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$   $b^0 = [0]$

and we take  $y_t / (\text{learning rate}) = 1$

After first iteration

$$w^1 = \begin{bmatrix} 124 \\ -113 \\ 216 \\ 5 \end{bmatrix} \quad b^1 = -6$$

After second iteration

$$w^2 = \begin{bmatrix} 123 \\ -105 \\ 216 \\ -8 \end{bmatrix} \quad b^2 = -7$$

After third iteration

$$w^3 = \begin{bmatrix} 121 \\ -106 \\ 218 \\ -17 \end{bmatrix} \quad b^3 = -7$$

Find iteration.

$$w = \begin{bmatrix} 111 \\ -106 \\ 217 \\ -50 \end{bmatrix} \quad b = -3$$

Number of iteration taken = 9.

Each iteration will includes summation of errors  
considering the complete data set and  
1 iteration to check the convergence

## ② Stochastic Gradient descent

In stochastic gradient descent, instead of taking the whole sum of indices we take a few indices ~~and~~ randomly and compute the sum and then minimize the function accordingly.

In this question we are taking data points one after 'another' and repeat accordingly. Therefore the equation used is

$$\begin{aligned} w^{(t+1)} &= w^{(t)} + \gamma_t y_i x_i \\ b^{(t+1)} &= b^{(t)} + \gamma_t y_i \end{aligned} \quad \left. \begin{array}{l} \text{if misclassified} \\ \text{otherwise} \end{array} \right\}$$

$$\begin{aligned} w^{(t+1)} &= w^t \\ b^{(t+1)} &= b^t \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

We start with  $w^0 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$        $b^0 = [0]$

and  $\gamma_t$  (learning rate) = 1

After first iteration

$$w^1 = \begin{bmatrix} -3 \\ 2 \\ 4 \\ 3 \end{bmatrix}, b^1 = -1$$

After second iteration

$$w^2 = \begin{bmatrix} 1 \\ 1 \\ 4 \\ 6 \end{bmatrix}, b^2 = -1$$

After third iteration

$$w^3 = \begin{bmatrix} 5 \\ -2 \\ 6 \\ 2 \end{bmatrix}, b^3 = -1$$

Final iteration

$$w = \begin{bmatrix} 7 \\ -4 \\ 11 \\ 0 \end{bmatrix}, b = 2$$

Number of iterations taken = 112

Each iteration is taken into consideration at every data point. Total iteration includes 100 iteration of complete classification.

### Problem 3

#### Support Vector Machines

In SVM we have to

$$\min_{\omega} \frac{1}{2} \|\omega^2\|$$

such that

$$y_i (\omega^T x^{(i)} + b) \geq 1, \text{ for all } i$$

∴ Using Lagrangian, we get

$$L(\omega, b, \lambda) = \frac{1}{2} \omega^T \omega + \sum_i \lambda_i (1 - y_i (\omega^T x^{(i)} + b))$$

$$\therefore \frac{\partial L}{\partial \omega_k} = \omega_k + \sum_i -\lambda_i y_i x_k^{(i)} = 0$$

$$\therefore \frac{\partial L}{\partial b} = \sum_i -\lambda_i y_i = 0$$

$$\therefore \omega = \sum_i \lambda_i y_i x^{(i)}$$

$$\sum_i \lambda_i y_i = 0$$

$\therefore$  substituting we get the dual as

$$\max_{\lambda \geq 0} -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^{(i)T} x_j + \sum_i \alpha_i$$

such that

$$\sum_i \alpha_i y_i = 0$$

Using the kernel trick, we have to take dot product now.

Data set:

	$x_1$	$x_2$	$x_3$	$y$
$s_1$	0	0	1	1
$s_2$	-1	-1	-1	1
$s_3$	1	1	-2	-1

We will use this vectors as they are assumed as support vectors.

Now we augment 1 as a bias input  $\therefore$

$$\therefore s_1 = (0 \ 0 \ 1 \ 1)$$

$$s_2 = (-1 \ -1 \ -1 \ 1)$$

$$s_3 = (1 \ 1 \ -2 \ 1)$$

$$\begin{aligned}\therefore \lambda_1 s_1 \cdot s_1 + \lambda_2 s_2 \cdot s_1 + \lambda_3 s_3 \cdot s_1 &= 1 \\ \lambda_1 s_1 \cdot s_2 + \lambda_2 s_2 \cdot s_2 + \lambda_3 s_3 \cdot s_2 &= 1 \\ \lambda_1 s_1 \cdot s_3 + \lambda_2 s_2 \cdot s_3 + \lambda_3 s_3 \cdot s_3 &= -1\end{aligned}$$

$$\Rightarrow \begin{aligned}2\lambda_1 + \lambda_2 \cdot 0 + \lambda_3 \cdot (-1) &= 1 \\ \lambda_1 \cdot 0 + \lambda_2 \cdot 4 + \lambda_3 \cdot 1 &= 1 \\ \lambda_1 \cdot (-1) + \lambda_2 \cdot 1 + \lambda_3 \cdot 7 &= -1\end{aligned}$$

$$\Rightarrow \begin{aligned}2\lambda_1 - \lambda_3 &= 1 \\ 4\lambda_2 + \lambda_3 &= 1 \\ -\lambda_1 + \lambda_2 + 7\lambda_3 &= -1\end{aligned}$$

$$\begin{aligned}\lambda_3 &= 2\lambda_1 - 1 \\ 4\lambda_2 + 2\lambda_1 - 1 &= 1 \Rightarrow 4\lambda_2 + 2\lambda_1 = 2 \\ &\Rightarrow 2\lambda_2 + \lambda_1 = 1 \\ &\Rightarrow \lambda_2 = \frac{1-\lambda_1}{2}\end{aligned}$$

$$\begin{aligned}-\lambda_1 + \lambda_2 + 7\lambda_3 &= -1 \\ -\lambda_1 + \frac{(1-\lambda_1)}{2} + 7(2\lambda_1 - 1) &= -1 \quad ] \times 2 \Rightarrow\end{aligned}$$

$$\begin{aligned}-2\lambda_1 + (1-\lambda_1) + 14(2\lambda_1 - 1) &= -2 \\ -2\lambda_1 + 1 - \lambda_1 + 28\lambda_1 - 14 &= -2\end{aligned}$$

$$25\lambda_1 = 11$$

$$\lambda_1 = \frac{11}{25} \quad \lambda_2 = \frac{14}{250} \quad \lambda_3 = -\frac{3}{25}$$

$$\text{As } \vec{w} = \sum_i \lambda_i \vec{s}$$

$$= \frac{11}{25} \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} + \frac{14}{50} \begin{pmatrix} -1 \\ -1 \\ 1 \end{pmatrix} + \frac{-3}{25} \begin{pmatrix} 1 \\ 1 \\ -2 \end{pmatrix}$$

$$= \begin{pmatrix} \frac{-14-6}{50} \\ \frac{-14-6}{50} \\ \frac{22-14+12}{50} \end{pmatrix} = \begin{pmatrix} -\frac{20}{50} \\ -\frac{20}{50} \\ \frac{20}{50} \end{pmatrix} = \begin{pmatrix} -\frac{2}{5} \\ -\frac{2}{5} \\ \frac{2}{5} \end{pmatrix}$$

As our vectors are augmented with a bias.

$$\therefore \vec{w} = \begin{bmatrix} -2/5 \\ -2/5 \\ 2/5 \end{bmatrix} \quad b = \begin{bmatrix} +3/5 \end{bmatrix}$$

the size of the margin is  $\frac{2}{\|\vec{w}\|}$ .

$$\|\vec{w}\| = \frac{2}{\sqrt{(2/5)^2 + (-2/5)^2 + (2/5)^2}}$$

$$\approx \frac{2}{\sqrt{3 \times 4/25}} = \frac{2}{2/5 \sqrt{3}}$$

$$= \frac{5}{\sqrt{3}} \approx 2.88$$

$\therefore$  the size of the margin is  $2.88 \approx 5/\sqrt{3}$

Problem 4:

## Lagrange Multipliers and Duality

Primal problem

$$\min_{p_1, \dots, p_n} \sum_{i=1}^n p_i \log p_i$$

such that  $p_1, \dots, p_n \geq 0$  and  $\sum_{i=1}^n p_i = 1$

Now we have to construct a dual optimization problem.

Let's first take the Lagrangian

$$L(p, v) = \sum_{i=1}^n p_i \log p_i + v \left[ \sum_{i=p}^n p_i - 1 \right]$$

$$\begin{aligned} \frac{\partial L}{\partial p_k} &= \log p_k + 1 + v = 0 \Rightarrow v = -1 - \log p_k \\ &\Rightarrow p_k = e^{-(1+v)} \end{aligned}$$

$$\frac{\partial L}{\partial v} = \sum_{i=1}^n p_i - 1 = 0 \Rightarrow \sum_{i=1}^n p_i = 1$$

$$\text{as } v = -1 - \log p_k$$

For  $k=1, 2, \dots, n \Rightarrow v$  will be same

$$\begin{aligned} \therefore p_1, p_2, p_3, \dots, p_n \text{ all have to be same} \\ \therefore p_1 = p_2 = p_3 = \dots = p_n = p \end{aligned}$$

As we also know

$$\sum_{i=1}^n p_i = 1 \Rightarrow p \cdot n = 1 \\ \Rightarrow p = \frac{1}{n}$$

$$\Rightarrow p_1 = p_2 = \dots = p_n = p = \frac{1}{n}$$

$$\therefore v = -1 - \log \frac{1}{n}$$

$$= -1 + \log n = (\log n + 1)$$

$$\Rightarrow v = (\log n + 1)$$

$$g(v) = \inf_p L(p, v)$$

$$= \sum_{i=1}^n \left[ e^{-(1+v)} \log e^{-1} \right] + e^{-(1+v)} \left[ \sum_{i=1}^n e^{-(1+v)} - 1 \right]$$

$$= n \left[ e^{-(1+v)} \cdot (-1)(1+v) \right] + e^{-(1+v)} \left[ n \cdot e^{-(1+v)} - 1 \right]$$

$$g(v) = -ne^{-(1+v)} - vne^{-(1+v)} + ne^{-(1+v)} - e^{-(1+v)}$$

To solve in  $g(v)$  substitute  $v = (\log n + 1)$

~~$$\therefore -ne^{-(1+\log n+1)} - (\log n + 1)ne^{-(1+v)}$$~~

First solve  $e^{-(1+v)} = e^{-(1+\log n+1)}$

$$\begin{aligned} &= e^{-(2+\log n)} \\ &= e^{-(\log(e^2 n))} \\ &= e^{(\log(e^2 n))} \\ &= (e^2 n)^{-1} \end{aligned}$$

~~$$\therefore -x \cdot e^{-2} - (log n + 1) \cdot e^{-2} x - x + e^{-2} \cdot n^{-2} \cdot n - e^{2} n^{-1}$$~~

$$\begin{aligned} &= -e^{-2} - e^{-2} \log n - e^{-2} + e^{-4} n^{-1} - e^2 n^{-1} \\ &= -e^{-2} [2 + \log n + e^{-2} n^{-1} + e^2 n^{-1}] \end{aligned}$$

Thus the solution of the primal problem  
is

$$-e^{-2} \left[ 2 + \log n - e^{-2} n^{-1} + n^{-2} \right]$$