

# Assignment 1

## Question 1:

The volume of a ball in  $p$  dimensional space is

$$V_p(R) = \frac{\pi^{p/2}}{\Gamma(\frac{p}{2}+1)} R^p$$

$\frac{\pi^{p/2}}{\Gamma(\frac{p}{2}+1)}$  part is constant given ball is  $p$  dimensional, Radius  $R$  changes.

We have  $N$  datapoints which are uniformly distributed in a unit ball. ( $R=1$ ). The probability that a point falls inside a ball of radius  $r$  can be found through the cdf:

$$\begin{aligned} P(X \in V_p(r)) &= \frac{V_p(r)}{V_p(1)} \\ &= \frac{\frac{\pi^{p/2}}{\Gamma(\frac{p}{2}+1)} r^p}{\frac{\pi^{p/2}}{\Gamma(\frac{p}{2}+1)} 1^p} \\ &= r^p \end{aligned}$$

Both the ball of radius  $r$  and unit ball has the same origin. So, if a point is inside a ball its distance from the origin is smaller than the radius of this ball.

Euclidean norm in  $p$  dimension space of vector  $x = \|x\|$

$$\text{So, } P(x \in V_p(r)) = P(\|x\| \leq r) = r^p$$

Now, let us consider a random variable

20

$$R = \min_i \{ \|x_i\| \}$$

which denotes smallest values of all samples.

$$P(R \leq r) = P(\min_i \{ \|x_i\| \} \leq r)$$

$$= 1 - P(\min_i \{ \|x_i\| \} > r)$$

$$= 1 - \prod_{i=1}^N P(\|x_i\| > r)$$

$$= 1 - \prod_{i=1}^N (1 - P(\|x_i\| \leq r))$$

$$= 1 - \prod_{i=1}^N (1 - r^p)$$

$$= 1 - (1 - r^p)^N$$

This cdf of  $R$  denotes the probability that the shortest distance of a point to origin in all random variable is less than or equal to cdf's parameter  $r$ .

Now median is 0.5 quantile of its cdf. So,

$$1 - (1 - r^p)^N = 0.5$$

$$\Rightarrow (1 - r^p)^N = \frac{1}{2}$$

$$\Rightarrow 1 - r^p = \frac{1}{2}^{\frac{1}{N}}$$

$$\Rightarrow r^p = 1 - \frac{1}{2}^{\frac{1}{N}}$$

$$\Rightarrow r = \left( 1 - \frac{1}{2}^{\frac{1}{N}} \right)^{\frac{1}{p}}$$

This, The median distance from origin to the closest data point is,

$$d(p, N) = \left(1 - \frac{1}{2} \frac{1}{N}\right)^{\frac{1}{p}}$$

For  $N = 10000$

$p = 1000$

$$d(1000, 10000) = \left(1 - 0.5 \frac{1}{10000}\right)^{\frac{1}{1000}} \quad (\text{proved})$$

Question 02:  $= \left(1 - 0.5 \frac{1}{10000}\right)^{\frac{1}{1000}} = \boxed{0.9905} \text{ (ans.)}$

Given,  $f(x) = (x_1 + x_2)(x_1 x_2 + x_1 x_2^2)$

$$= x_1^2 x_2 + x_1^2 x_2^2 + x_1 x_2^2 + x_1 x_2^3$$

Now, the gradient,  $\nabla f(x) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x) \\ \frac{\partial}{\partial x_2} f(x) \end{bmatrix}$

$$\therefore \nabla f(x) = \begin{bmatrix} 2x_1 x_2 + 2x_1 x_2^2 + x_2^2 + x_2^3 \\ x_1^2 + 2x_1 x_2 + 2x_1 x_2^2 + 3x_1 x_2^2 \end{bmatrix}$$

Finding 3 stationary points,  $\nabla f(x) = 0$

$$\therefore 2x_1 x_2 + 2x_1 x_2^2 + x_2^2 + x_2^3 = 0$$

$$\Rightarrow 2x_1 (x_2 + x_2^2) + x_2 (x_2 + x_2^2) = 0$$

$$\Rightarrow (x_2 + x_2^2) (2x_1 + x_2) = 0$$

$$2x_1 + x_2 = 0 \quad \left| \quad x_2 + x_2^2 = 0 \right.$$

$$x_2 = -2x_1$$

$$x_2 (1 + x_2) = 0$$

$$x_2 = 0 \quad x_2 = -1$$

$$x_1^2 + 2x_1^2 x_2 + 2x_1 x_2 + 3x_1 x_2^2 = 0$$

$$x_1 (x_1 + 2x_1 x_2 + 2x_2 + 3x_2^2) = 0$$

or

$$x_1 = 0$$

$$x_1 + 2x_1 x_2 + 2x_2 + 3x_2^2 = 0 \quad \text{--- (1)}$$

From previous page,

$$\text{let, } x_2 = 0 \text{ so from (1) } x_1 = 0$$

$$\text{let } x_2 = -1 \text{ so from (1) } x_1 - 2x_1 - 2 + 3 = 0$$

$$\Rightarrow -x_1 + 1 = 0$$

$$\therefore x_1 = 1$$

$$\text{Let } x_2 = -2x_1 \text{ so from (1)}$$

$$x_1 - 4x_1^2 - 4x_1 + 12x_1^2 = 0$$

$$\Rightarrow -3x_1 + 8x_1^2 = 0$$

$$\Rightarrow x_1 (-3 + 8x_1) = 0$$

$$x_1 = 0$$

$$-3 + 8x_1 = 0$$

$$x_1 = \frac{3}{8}$$

$$x_2 = -2 \cdot x_1 = -2 \cdot \frac{3}{8} = -\frac{6}{8}$$

So, stationary points,  $(x_1, x_2) = (0, 0), (1, -1), \left(\frac{3}{8}, -\frac{6}{8}\right)$



Now, the Hessian matrix:

$$\nabla^2 f(x) = \begin{bmatrix} 2x_2 + 2x_2^2 & 2x_1 + 4x_1x_2 + 2x_2 + 3x_2^2 \\ 2x_1 + 4x_1x_2 + 2x_2 + 3x_2^2 & 2x_1^2 + 2x_1 + 6x_1x_2 \end{bmatrix}$$

at point  $(0, 0)$

$$\nabla^2 f(x) = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

If a Hessian matrix is negative definite then the corresponding critical points are local maxima.

We can check if a matrix is negative definite through its eigenvalues.

The eigen value of the Hessian matrix at  $0, 0$  is

$[0, 0]$  so it's not negative definite.

at point  $(1, -1)$

$$\nabla^2 f(x) = \begin{bmatrix} -2 + 2 & 2 + 4 \cdot 1 \cdot (-1) - 2 + 3 \\ 2 - 4 - 2 + 3 & 2 + 2 - 6 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & -1 \\ -1 & -2 \end{bmatrix}$$

Its eigenvalue  $[0.414 \quad -2.414]$

So it's also not a local maxima. Rather a saddle point.

Now, at point  $\left(\frac{3}{8}, -\frac{6}{8}\right)$

$$\nabla^2 f(x) = \begin{bmatrix} -\frac{2 \times 6}{8} + 2 \cdot \left(\frac{6}{8}\right)^2 & 2 \cdot \frac{3}{8} + 4 \cdot \frac{3}{8} \left(-\frac{6}{8}\right) + 2 \left(-\frac{6}{8}\right) + 3 \left(-\frac{6}{8}\right) \\ 2 \cdot \frac{3}{8} + 4 \cdot \frac{3}{8} \left(-\frac{6}{8}\right) + 2 \left(-\frac{6}{8}\right) + 3 \left(-\frac{6}{8}\right) & -\frac{21}{32} \end{bmatrix}$$

$$= \begin{bmatrix} -\frac{3}{8} & -\frac{3}{16} \\ -\frac{3}{16} & -\frac{21}{32} \end{bmatrix}$$

Its corresponding eigenvalue is:  $[-2.8125, -0.75]$   
 So the Hessian is negative definite.

So the point  $\left(\frac{3}{8}, -\frac{6}{8}\right)$  is local maximum.

and it is the only local maximum of the given function.

Question 08:

$$f(x) = 8x_1 + 12x_2 + x_1^2 - 2x_2^2$$

$$\nabla f(x) = \begin{bmatrix} 8 + 0 + 2x_1 + 0 \\ 0 + 12 + 0 - 4x_2 \end{bmatrix} = \begin{bmatrix} 2x_1 + 8 \\ -4x_2 + 12 \end{bmatrix}$$

$$\nabla^2 f(x) = \begin{bmatrix} 2 & 0 \\ 0 & -4 \end{bmatrix}$$

Now there is only constants in the Hessian which implies there is only one solution to the given quadratic equation.

The stationary point

$$\begin{aligned} \nabla f(x) = 0 \quad 2x_1 + 8 = 0 &\Rightarrow x_1 = -4 \\ -4x_2 + 12 = 0 &\Rightarrow x_2 = 3 \end{aligned}$$

The stationary point doesn't have any effect on the Hessian.

The eigenvalue of the Hessian is  $\begin{bmatrix} 2 & -4 \end{bmatrix}$

There is a positive eigenvalue and also a ~~mi~~ negative eigenvalue. So, the Hessian is neither positive or negative definite.

So, the stationary point is a saddle point.

(proved).

### Question 04:

Let  $A$  be a  $n \times n$  matrix

$B$  be a  $m \times m$  matrix.

$$A = \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ A_{21} & \dots & \dots & A_{2n} \\ \vdots & & & \\ A_{n1} & \dots & \dots & A_{nn} \end{bmatrix}$$

$$B = \begin{bmatrix} B_{11} & \dots & B_{1m} \\ B_{21} & \dots & B_{2m} \\ \vdots & & \vdots \\ B_{m1} & \dots & B_{mm} \end{bmatrix}$$

as  $A$  is positive definite we can say

$$x'^T A x' > 0 \quad \text{where } x' \text{ is its corresponding eigenvector}$$

Similarly  $B$  is also ~~semi~~ positive definite

$$x''^T B x'' > 0 \quad \text{where } x'' \text{ is } B\text{'s eigenvector}$$

$$x' = \begin{bmatrix} x'_1 \\ x'_2 \\ \vdots \\ x'_n \end{bmatrix} \quad x'' = \begin{bmatrix} x''_1 \\ x''_2 \\ \vdots \\ x''_m \end{bmatrix}$$

We can represent  $x'^T A x' > 0$  in terms of matrix multiplied form as the following

$$\sum_{j=1}^n \sum_{i=1}^n x'_i A_{ij} x'_j > 0 \quad \text{--- (1)}$$



similarly  $x'^T B x'' > 0$  as

$$\sum_{j=1}^m \sum_{i=1}^m x''_i B_{ij} x''_j > 0 \quad \text{--- (11)}$$

Now let's consider  $C = \begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix}$  which can be expressed as the following:

$$C = \begin{bmatrix} \underbrace{\begin{matrix} A_{11} & \dots & A_{1n} \\ A_{21} & \dots & A_{2n} \\ \vdots & & \vdots \\ A_{n1} & \dots & A_{nn} \end{matrix}}_{n \times n} & \underbrace{\begin{matrix} 0 & \dots & 0 \\ 0 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & 0 \end{matrix}}_{n \times m} \\ \underbrace{\begin{matrix} 0 & \dots & 0 \\ 0 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & 0 \end{matrix}}_{m \times n} & \underbrace{\begin{matrix} B_{11} & \dots & B_{1m} \\ B_{21} & \dots & B_{2m} \\ \vdots & & \vdots \\ B_{m1} & \dots & B_{mm} \end{matrix}}_{m \times m} \end{bmatrix}$$

All 0's      All 0's

dimension

$$(n+m) \times (n+m)$$

We want to prove  $C$  is also positive definite. So,  $C$  also has an eigenvector  $x$  for which  $x^T C x > 0$ .

where,

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \\ x_{n+1} \\ \vdots \\ x_{n+m} \end{bmatrix} = \begin{bmatrix} x'_1 \\ x'_2 \\ \vdots \\ x'_n \\ x''_1 \\ x''_2 \\ \vdots \\ x''_m \end{bmatrix} = \begin{bmatrix} x' \\ x'' \end{bmatrix}$$

So,  $x$  is the join of both eigenvectors of  $A$  and  $B$ .

If  $x^T C x > 0$  then we would be able to say  $C$  is positive definite.

$$x^T C x$$

$$[x_1 \ x_2 \ \dots \ x_n \ x_{n+1} \ \dots \ x_{n+m}]$$

$$C = \begin{bmatrix} \boxed{A} & \boxed{0} \\ \boxed{0} & \boxed{B} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \\ x_{n+1} \\ \vdots \\ x_{n+m} \end{bmatrix}$$

$\begin{matrix} & \begin{matrix} n & n+1 & n+m \end{matrix} \\ \begin{matrix} \uparrow \\ n \\ n+1 \\ n+m \end{matrix} & \end{matrix}$

doing matrix multiplication we can get the following result:

$$\sum_{j=1}^n \left[ \sum_{i=1}^n x_i C_{ij} x_j + \sum_{i=n+1}^{n+m} x_i \cdot 0 \cdot x_j \right]$$

$$+ \sum_{j=n+1}^{n+m} \left[ \sum_{i=1}^n x_i \cdot 0 \cdot x_j + \sum_{i=n+1}^{n+m} x_i \cdot C_{ij} \cdot x_j \right]$$

$$\Rightarrow \underbrace{\sum_{j=1}^n \sum_{i=1}^n x_i C_{ij} x_j}_{\text{This part of } C \text{ corresponds to } A \text{ and is equal to (I)}} + \underbrace{\sum_{j=n+1}^{n+m} \sum_{i=n+1}^{n+m} x_i C_{ij} x_j}_{\text{This part of } C \text{ corresponds to } B \text{ and is equal to (II)}}$$

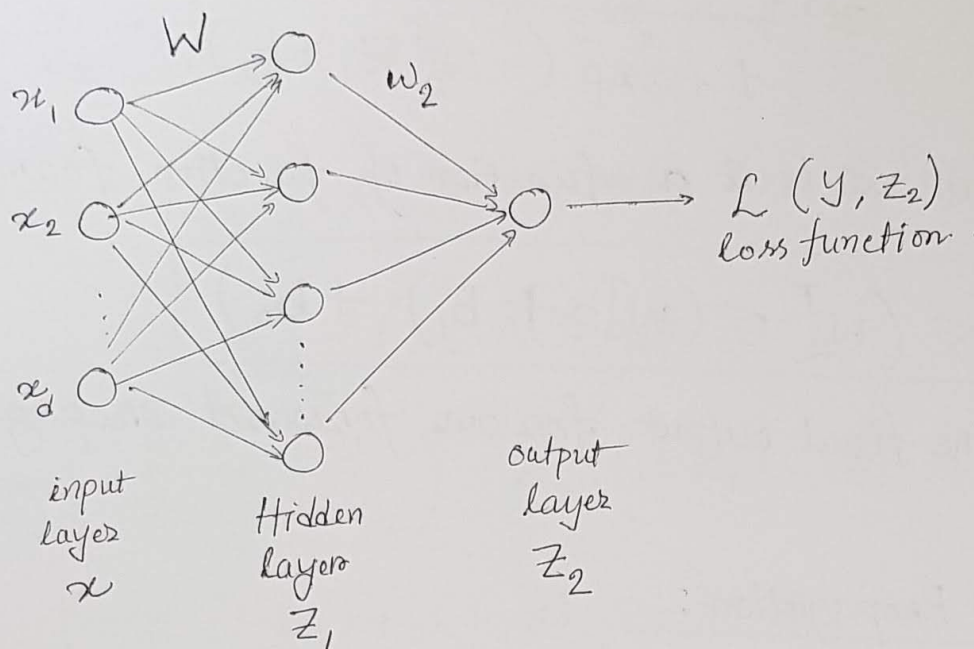
This part of  $C$  corresponds to  $A$  and is equal to (I)

This part of  $C$  corresponds to  $B$  and is equal to (II)

as (I) and (II) are  $> 0$  so addition of (I) + (II)  $> 0$ .

$$\text{So, } x^T C x > 0$$

$\therefore C$  or  $\begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix}$  is a positive definite matrix (proved)

Forward Propagation:

The input is a vector  $x$ .

$i^{\text{th}}$  node of  $z_1^{[i]}$  is calculated the following way:

[we are using sigmoid function for activation and  $b_1$  is bias term in 1<sup>st</sup> layer along with parameter  $W$  where  $W^{[i]}$  is parameter corresponding to  $i$ -th node in  $z_1^{[i]}$ ]

$$z_1^{[i]} = \sigma(W^{[i]T}x + b_1) = \frac{1}{1 + \exp(-W^{[i]T}x + b_1)}$$

Vectorizing for the entire layer we get:

$$z_1 = \sigma(W^T x + b_1) = \frac{1}{1 + \exp(-W^T x + b_1)}$$

Now for the output layer  $z_1$  is used as input and  $w_2$ ,  $b_2$  is our parameter. Similarly using sigmoid



We get  $z_2 = \sigma(w_2^T z_1 + b_2)$

$$= \frac{1}{1 + \exp(-w_2^T z_1 + b_2)}$$

OR we can write it as function of function form in below:

$$\boxed{z_2 = \sigma(w_2^T \sigma(w^T x + b_1) + b_2)}$$

Which is the final output for our forward propagation.

Backward Propagation:

Our loss function is a 2-class cross-entropy loss function

$$L(z_2, y) = -(y \log z_2 + (1-y) \log(1-z_2))$$

Now we start backprop from our last layer, so

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}$$

Using chain rule.

$$\begin{aligned} \frac{\partial L}{\partial z_2} &= \frac{\partial}{\partial z_2} (y \log z_2 + (1-y) \log(1-z_2)) \\ &= -\frac{y}{z_2} + \frac{1-y}{1-z_2} \end{aligned}$$



$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$

$$= \left( -\frac{y}{z_2} + \frac{1-y}{1-z_2} \right) \frac{\partial}{\partial w_2} \sigma(w_2^T z_1 + b_2)$$

$$= \left( -\frac{y}{z_2} + \frac{1-y}{1-z_2} \right) \sigma(w_2^T z_1 + b_2) (1 - \sigma(w_2^T z_1 + b_2))$$

$$\cdot \frac{\partial}{\partial w_2} (w_2^T z_1 + b_2)$$

$$= \left( -\frac{y}{z_2} + \frac{1-y}{1-z_2} \right) \cdot z_2 (1-z_2) \cdot z_1$$

$$= (-y(1-z_2) + (1-y) \cdot z_2) z_1$$

$$= (-y + y \cdot z_2 + z_2 - y \cdot z_2) \cdot z_1$$

$$= (z_2 - y) \cdot z_1$$

To match dimension of  $\frac{\partial \mathcal{L}}{\partial w_2}$  with  $w_2$  we rearrange and get,

$$\boxed{\frac{\partial \mathcal{L}}{\partial w_2} = z_1(z_2 - y)}$$

For parameter  $b_2$

$$\frac{\partial \mathcal{L}}{\partial b_2} = \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial b_2}$$

$$\frac{\partial z_2}{\partial b_2} = \frac{\partial}{\partial b_2} \sigma(w_2^T z_1 + b_2) = z_2 (1 - z_2)$$

$$\boxed{\frac{\partial \mathcal{L}}{\partial b_2} = (z_2 - y)}$$

Now for the parameters of the 1<sup>st</sup> hidden layer:

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z_2} \cdot \frac{\partial z_2}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

$$\begin{aligned} \frac{\partial}{\partial z_1} z_2 &= \frac{\partial}{\partial z_1} \sigma(w_2^T z_1 + b_2) \\ &= \sigma(w_2^T z_1 + b_2) (1 - \sigma(w_2^T z_1 + b_2)) \frac{\partial}{\partial z_1} (w_2^T z_1 + b_2) \\ &= z_2 (1 - z_2) \cdot w_2 \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial w_1} z_1 &= \frac{\partial}{\partial w} \sigma(w^T x + b_1) \\ &= z_1 (1 - z_1) \cdot \frac{\partial}{\partial w} (w^T x + b_1) \\ &= z_1 (1 - z_1) \cdot x \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w_1} &= \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial w} \\ &= \left( -\frac{y}{z_2} + \frac{1-y}{1-z_2} \right) z_2 (1-z_2) w_2 z_1 (1-z_1) x \end{aligned}$$

$$\boxed{\frac{\partial \mathcal{L}}{\partial w} = z_1 (z_2 - y) \cdot w_2 \cdot z_1 (1 - z_1) \cdot x}$$

To match the dimension of  $\frac{\partial \mathcal{L}}{\partial w}$  with  $w$  we rearrange and get the following.

$$\frac{\partial \mathcal{L}}{\partial w} = x [z_1 (z_2 - y)]^T [z_1 \odot (1 - z_1)] w_2^T$$

For parameter  $b_1$

$$\frac{\partial \mathcal{L}}{\partial b_1} = \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

$$\frac{\partial}{\partial b_1} z_1 = \frac{\partial}{\partial b_1} \sigma(w^T x + b_1)$$

$$= z_1 (1 - z_1) \frac{\partial}{\partial b_1} (w^T x + b_1)$$

$$= z_1 (1 - z_1)$$

$$\therefore \boxed{\frac{\partial \mathcal{L}}{\partial b_1} = z_1 (z_2 - y) \cdot w_2 \cdot z_1 (1 - z_1)}$$

All gradients based on all parameters  $w_1, b_1, w_2, b_2$  has been shown.