

Assignment 7

ECSE 4200

{ Name: Fikero Lin
NetID: f1955 }

Problem #1:

- $\vec{x} \in \mathbb{R}^d$; all entries non-zero
- $w \in \mathbb{R}^d$
- $\vec{y} = \max \{0, w\vec{x}\}$
- w follows $N(0, 1)$ distribution.

1. $\vec{y} = \max \{0, w\vec{x}\}$

$$\vec{x} = \left[\begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_d \end{array} \right] \Bigg\}^d$$

$$w\vec{x} = \begin{bmatrix} w_{11} & \dots & w_{1d} \\ \vdots & & \vdots \\ w_{m1} & \dots & w_{md} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$w = \underbrace{\begin{bmatrix} w_{11} & \dots & w_{1d} \\ \vdots & & \vdots \\ w_{m1} & \dots & w_{md} \end{bmatrix}}_d \Bigg\}^m$$

$$= \begin{bmatrix} x_1 w_{11} + \dots + x_d w_{1d} \\ \vdots \\ x_1 w_{m1} + \dots + x_d w_{md} \end{bmatrix}$$

$$\hat{y} = \max \{0, \tilde{w}^T x\}$$

$$= \max \left\{ 0, \begin{bmatrix} x_1 w_{11} + \dots + x_d w_{1d} \\ \vdots \\ x_1 w_{m1} + \dots + x_d w_{md} \end{bmatrix} \right\}$$

Component-wise max function

Each entry of W is independent, distributed according to $N(0, 1)$:

$$\text{Let } \vec{w}_i = \begin{bmatrix} w_{i1} \\ w_{i2} \\ \vdots \\ w_{id} \end{bmatrix} = \begin{bmatrix} x_1 w_{i1} + \dots + x_d w_{id} \\ \vdots \\ x_1 w_{mi1} + \dots + x_d w_{mid} \end{bmatrix}$$

Since elements of W are normally distributed and independent:

$$\begin{aligned} \mathbb{E}(a_1) &= \mathbb{E}(x_1 w_{11} + \dots + x_d w_{1d}) \quad \text{linearity} \\ &= \mathbb{E}(x_1 w_{11}) + \dots + \mathbb{E}(x_d w_{1d}) \\ &= \underbrace{x_1 \mathbb{E}(w_{11})}_{\text{mean is 0}} + \dots + \underbrace{x_d \mathbb{E}(w_{1d})}_{\text{mean is 0}} \end{aligned}$$

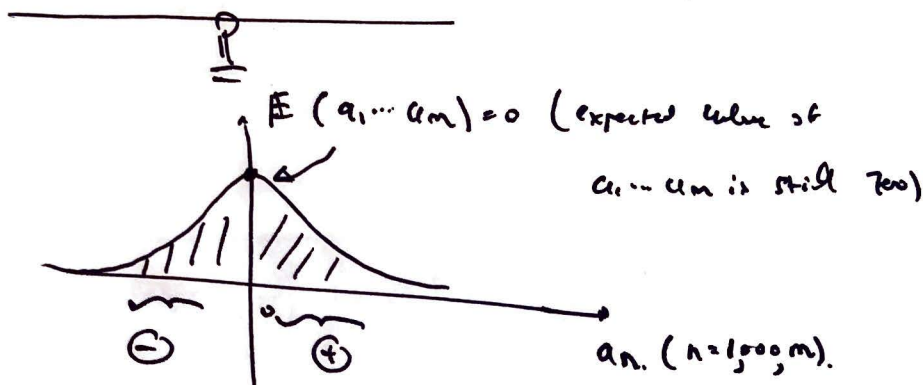
$$\boxed{\neq 0}$$

Similarly, $\mathbb{E}(a_2) \dots \mathbb{E}(a_m) = 0$

Thus, we see that a_1, \dots, a_m are again distributed

with mean of 0. In other words, we expected that

half of a_1, \dots, a_m are greater than 0 and half of a_1, \dots, a_m are less than 0.



Since $\bar{y} = \max\{0, \bar{x}\}$ and half of a_1, \dots, a_m is ^{expected to be} greater than 0 and half is expected to be less than 0 around half of entries in \bar{y} are 0.

\therefore Thus, $m/2$ entries in \bar{y} are expected to be non-zero.

$$2. \|x\|_2^2 = \sigma^2 = x_1^2 + \dots + x_d^2$$

$$w\bar{x} = \begin{bmatrix} a_1 \\ \vdots \\ a_m \end{bmatrix} = \begin{bmatrix} x_1 w_{11} + \dots + x_d w_{d1} \\ \vdots \\ x_1 w_{m1} + \dots + x_d w_{md} \end{bmatrix}$$

\Rightarrow the mean of a_1, \dots, a_m will be zero as shown before.

Note that:

$$\begin{aligned}
 & \mathbb{E}(c \cdot x) = c \cdot \mathbb{E}(x) \\
 & \text{var}(c \cdot x) = c^2 \cdot \text{var}(x) \\
 & \sigma^2(c \cdot x) = c^2 \sigma^2(x)
 \end{aligned}$$

$$\begin{aligned}
 & x \sim N(\mu_x, \sigma_x^2) \\
 & y \sim N(\mu_y, \sigma_y^2) \\
 & z = x + y \rightarrow z \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)
 \end{aligned}$$

$$\mathbf{z} = x_1 \mathbf{u}_1 + \dots + x_d \mathbf{u}_d$$

$$x_1 \mathbf{u}_1 \sim N(0, \sigma_1^2 \cdot 1) = \boxed{N(0, \sigma_1^2)}$$

$$x_2 \mathbf{u}_2 \sim N(0, \sigma_2^2 \cdot 1) = \boxed{N(0, \sigma_2^2)}$$

\vdots

$$x_d \mathbf{u}_d \sim \boxed{N(0, \sigma_d^2)}$$

$$\|\mathbf{x}\|^2 = \sigma^2$$

$$\text{Thus, } \mathbf{z} \sim N(0, \sigma_1^2 + \dots + \sigma_d^2)$$

$$\boxed{N(0, \sigma^2)}$$

\therefore the entries of \mathbf{w} are distributed according to: $N(0, \sigma^2)$.

3. Each entry in \vec{y} :

$$\vec{y} = \begin{bmatrix} \max(0, a_1) \\ \max(0, a_2) \\ \vdots \\ \max(0, a_m) \end{bmatrix} = \begin{bmatrix} a_1' \\ a_2' \\ \vdots \\ a_m' \end{bmatrix}$$

$a_i \sim \text{mean } N(0, \sigma^2)$

$p(\text{probability}) \left(\frac{e^{-a_i^2/2\sigma^2}}{\sqrt{2\pi}\sigma} \right)$

$$E(a_i') = \int_{-\infty}^{\infty} \max(0, a_i) p(a_i) da_i$$

< 0 part doesn't contribute anything

$$= \int_0^{\infty} \max(a_i) p(a_i) da_i$$

$$= \int_0^{\infty} a_i \cdot \frac{e^{-a_i^2/2\sigma^2}}{\sqrt{2\pi}\sigma} da_i$$

$$\boxed{\frac{\sigma}{\sqrt{2\pi}}}$$

\therefore the expected value / mean of each entry in \vec{y} is $\frac{\sigma}{\sqrt{2\pi}}$.

$$a_i \sim N(0, \sigma^2) \Rightarrow \frac{e^{-a_i^2/2\sigma^2}}{\sqrt{2\pi}\sigma}$$

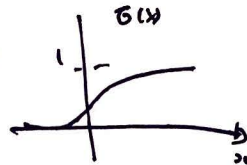
Problem 2:

$m=2, d=2.$

$w = \begin{bmatrix} 1 & 2 \\ -2 & 3 \end{bmatrix}, \vec{z} = \begin{bmatrix} 2 \\ -3 \end{bmatrix}$

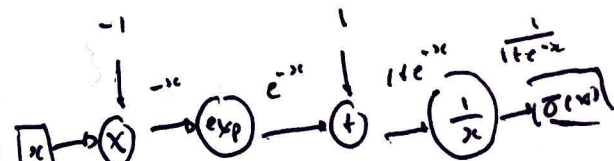
Let $L = \max \{ \sigma(w_{(1)} \vec{z}), \sigma(w_{(2)} \vec{z}) \}$

$\sigma(x) = \frac{1}{1 + e^{-x}}$

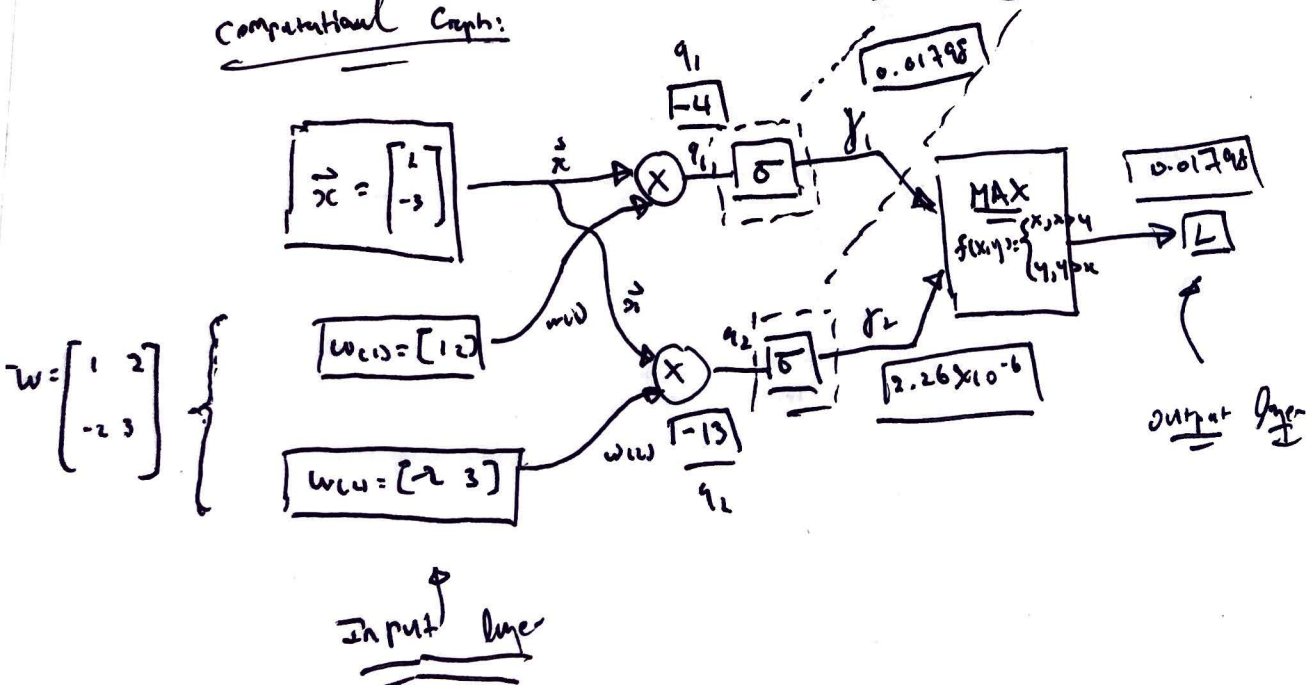


Wk0:

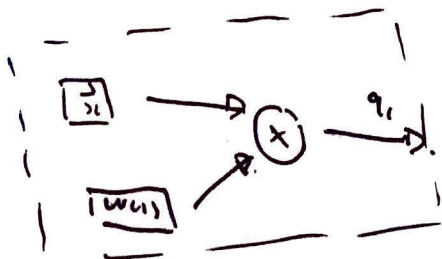
Row vector: $\begin{cases} w_{(1)} = [1 \ 2] = \begin{bmatrix} 1 \\ 2 \end{bmatrix}^T \\ w_{(2)} = [-2 \ 3] = \begin{bmatrix} -2 \\ 3 \end{bmatrix}^T \end{cases}$



Computational Graph:



Computation of gradients:



$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$w_{(1)} = [w_1^{(1)} \ w_2^{(1)}]$$

$$q_1 = w_{(1)} \cdot \vec{x}$$

$$q_1 = [w_1^{(1)} \ w_2^{(1)}] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = q_1(\vec{x}, w_{(1)})$$

$$= [x_1 w_1^{(1)} + x_2 w_2^{(1)}] = q_1$$

Jacobian:

$$\Rightarrow \left[\frac{\partial q_{(1)}}{\partial x_j} \right] = w_j^{(1)} \quad \left\{ \begin{array}{l} 1 \times 1 \rightarrow \text{treat as scalar} \\ \text{matrix} \end{array} \right.$$

$$w_1^{(1)} = \frac{\partial q_{(1)}}{\partial x_1}$$

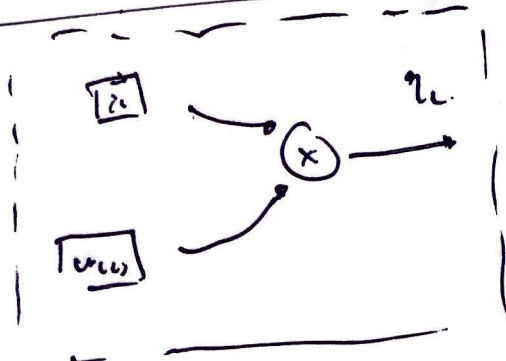
$$w_2^{(1)} = \frac{\partial q_{(1)}}{\partial x_2}$$

2nd row

$$\begin{array}{l} w_1^{(1)} = 1 = \frac{\partial q_{(1)}}{\partial x_1} \\ w_2^{(1)} = 2 = \frac{\partial q_{(1)}}{\partial x_2} \end{array}$$

$$\Rightarrow \left[\frac{\partial q_{(1)}(\vec{x}, w_{(1)})}{\partial w_{j(1)}} \right] = x_j \quad \left\{ \begin{array}{l} x_1 = \frac{\partial q_{(1)}}{\partial w_1^{(1)}} \\ x_2 = \frac{\partial q_{(1)}}{\partial w_2^{(1)}} \end{array} \right.$$

$$\begin{array}{l} x_1 = 2 = \frac{\partial q_{(1)}}{\partial w_1^{(1)}} \\ x_2 = -3 = \frac{\partial q_{(1)}}{\partial w_2^{(1)}} \end{array}$$



$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$w_{(2)} = [w_1^{(2)} \ w_2^{(2)}]$$

$$q_2 = [w_1^{(2)} \ w_2^{(2)}] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$= [x_1 w_1^{(2)} + x_2 w_2^{(2)}]$$

Jacobian:

$$\frac{\partial q_{2,1}(\bar{x}, w_{2,1})}{\partial x_j} = w_{j,2}$$

$$w_{1,2} = \frac{\partial q_{2,1}}{\partial x_1}$$

$$w_{2,2} = \frac{\partial q_{2,1}}{\partial x_2}$$

$$w_{1,2} = -2 = \frac{\partial q_{2,1}}{\partial x_1}$$

$$\therefore w_{2,2} = 3 = \frac{\partial q_{2,1}}{\partial x_2}$$

$$\frac{\partial q_{2,1}(\bar{x}, w_{2,1})}{\partial w_{j,2}} = x_j$$

$$x_1 = \frac{\partial q_{2,1}}{\partial w_{1,2}}$$

$$x_2 = \frac{\partial q_{2,1}}{\partial w_{2,2}}$$

$$x_1 = 2 = \frac{\partial q_{2,1}}{\partial w_{1,2}}$$

$$x_2 = -3 = \frac{\partial q_{2,1}}{\partial w_{2,2}}$$



$$\delta_1 = \sigma(q_1) = \frac{1}{1 + e^{-q_1}}$$

$$q_1 = -4$$

$$\frac{\partial \delta_1}{\partial q_1} = \frac{\partial}{\partial q_1} [\sigma(q_1)] = \sigma(q_1)(1 - \sigma(q_1))$$

$$\therefore \frac{\partial \delta_1}{\partial q_1} = \sigma(-4)(1 - \sigma(-4))$$

$$= 0.01766$$



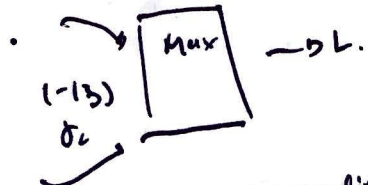
$$q_2 = -13$$

$$\frac{\partial \delta_2}{\partial q_2} = \sigma(q_2)(1 - \sigma(q_2))$$

$$\therefore \frac{\partial \delta_2}{\partial q_2} = \sigma(-13)(1 - \sigma(-13))$$

$$= 2.26 \times 10^{-6}$$

$$\delta_1 (-4)$$



δ_1 gradient passes if $\delta_1 > \delta_2$
 δ_2 gradient passes if $\delta_2 > \delta_1$

Since $\delta_1 > \delta_2$,

gradient of 0.01766
 passes through
 max gate.

problem #3:

$$\hat{y}(z_1, z_2) = \frac{e^{z_1}}{e^{z_1} + e^{z_2}}$$

Cross-entropy loss:

$$L(y, \hat{y}) = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

Goal: Show that $\frac{\partial L(y, \hat{y})}{\partial z_1} = \hat{y} - y$, $\frac{\partial L(y, \hat{y})}{\partial z_2} = y - \hat{y}$

$$\textcircled{1} \quad \frac{\partial L}{\partial z_1} = \frac{\partial}{\partial z_1} \left[-y \log(\hat{y}) - (1-y) \log(1-\hat{y}) \right]$$

y is not function of z_1 but \hat{y} is

$$= -y \frac{\partial}{\partial z_1} \log(\hat{y}) - (1-y) \frac{\partial}{\partial z_1} \log(1-\hat{y})$$

$$= -y \frac{\partial}{\partial z_1} \left[\log \left[\frac{e^{z_1}}{e^{z_1} + e^{z_2}} \right] \right] - (1-y) \frac{\partial}{\partial z_1} \left[\log \left(\frac{e^{z_2}}{e^{z_1} + e^{z_2}} \right) \right]$$

$$= -y \frac{\partial}{\partial z_1} \left[\log e^{z_1} - \log(e^{z_1} + e^{z_2}) \right] - (1-y) \frac{\partial}{\partial z_1} \left[\log e^{z_2} - \log(e^{z_1} + e^{z_2}) \right]$$

$$= -y \frac{\partial}{\partial z_1} \left[z_1 - \log(e^{z_1} + e^{z_2}) \right] - (1-y) \frac{\partial}{\partial z_1} \left[z_2 - \log(e^{z_1} + e^{z_2}) \right]$$

$$= -y \left[1 - \frac{e^{\hat{y}}}{e^{\hat{y}_1} + e^{\hat{y}_2}} \right] - (1-y) \left[0 - \frac{e^{\hat{y}_1}}{e^{\hat{y}_1} + e^{\hat{y}_2}} \right]$$

$$= -y + y \cdot \hat{y} + (1-y) \hat{y}$$

$$= -y + y \cdot \hat{y} + \hat{y} - y \hat{y} = \boxed{\hat{y} - y} \quad \therefore \text{Ans}$$

$$\textcircled{2}: \frac{\partial L}{\partial \hat{y}_1} = \frac{\partial}{\partial \hat{y}_1} \left[-y \log(\hat{y}) - (1-y) \log(1-\hat{y}) \right]$$

$$= -y \frac{\partial}{\partial \hat{y}_1} \log(\hat{y}) - (1-y) \frac{\partial}{\partial \hat{y}_1} \log(1-\hat{y})$$

$$= -y \left[\frac{\partial}{\partial \hat{y}_1} \left(\frac{e^{\hat{y}_1}}{e^{\hat{y}_1} + e^{\hat{y}_2}} \right) \right] - (1-y) \frac{\partial}{\partial \hat{y}_1} \left[\frac{e^{\hat{y}_2}}{e^{\hat{y}_1} + e^{\hat{y}_2}} \right]$$

$$= -y \frac{\partial}{\partial \hat{y}_1} \left[\log e^{\hat{y}_1} - \log(e^{\hat{y}_1} + e^{\hat{y}_2}) \right] - (1-y) \frac{\partial}{\partial \hat{y}_1} \left[\log e^{\hat{y}_2} - \log(e^{\hat{y}_1} + e^{\hat{y}_2}) \right]$$

$$= -y \left[0 - \frac{e^{\hat{y}_1}}{e^{\hat{y}_1} + e^{\hat{y}_2}} \right] - (1-y) \left[0 - \frac{e^{\hat{y}_2}}{e^{\hat{y}_1} + e^{\hat{y}_2}} \right]$$

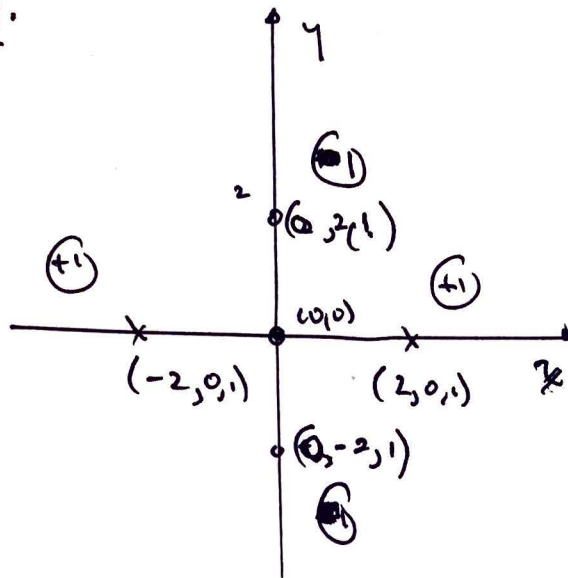
$$= y \frac{e^{\hat{y}_1}}{e^{\hat{y}_1} + e^{\hat{y}_2}} - (1-y) \left(1 - \frac{e^{\hat{y}_2}}{e^{\hat{y}_1} + e^{\hat{y}_2}} \right)$$

$$= y(1-\hat{y}) - (1-y)\hat{y}$$

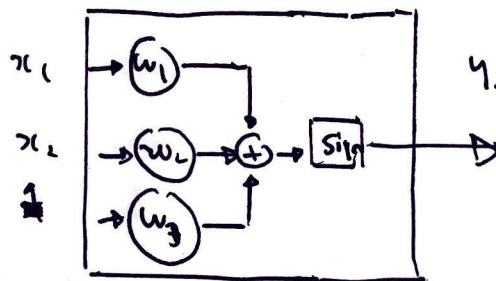
$$= y - y\hat{y} - \hat{y} + y\hat{y} = \boxed{y - \hat{y}} \quad \therefore \text{Ans}$$

Problem #4.

1.



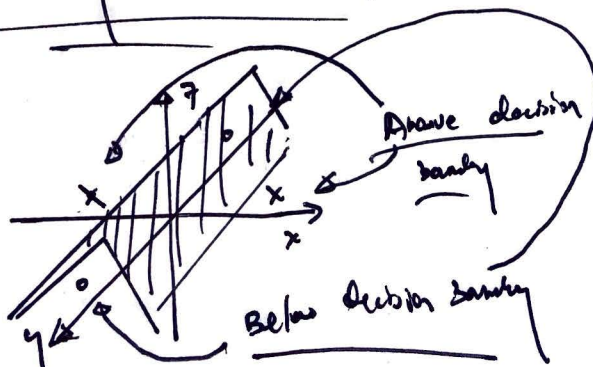
Unit:



$$y = \text{sign}(x_1 w_1 + x_2 w_2 + x_3 w_3)$$

The goal is to correctly classify all four points.
we can do this by creating decision boundaries in \mathbb{R}^3 space

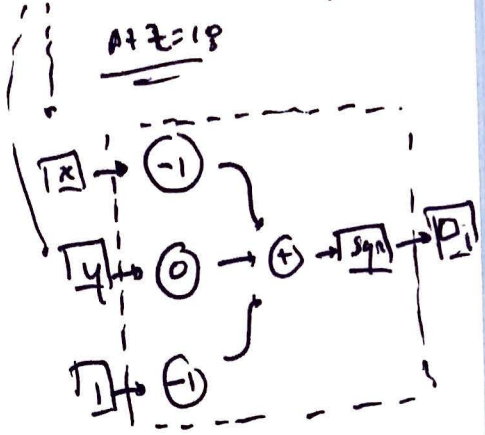
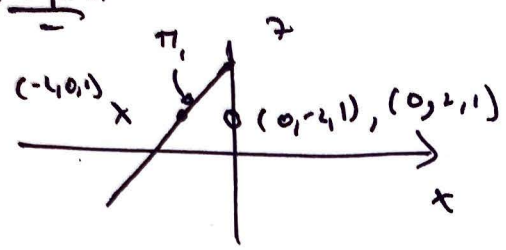
we have the following decision boundary:



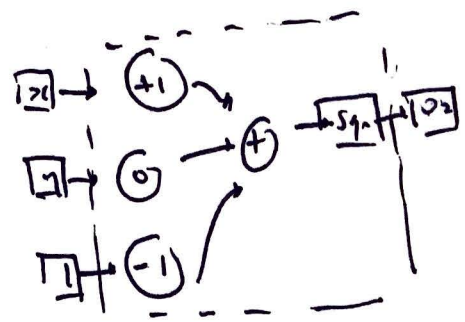
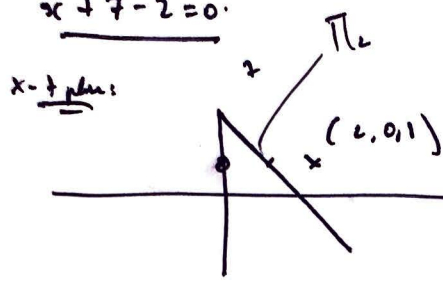
This decision boundary is composed of two planes:

(x,y) coord of points to classify

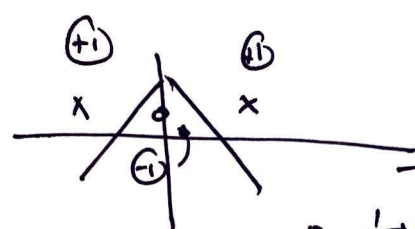
Plane π_1 : $-x + z - 2 = 0$
x-z plane:



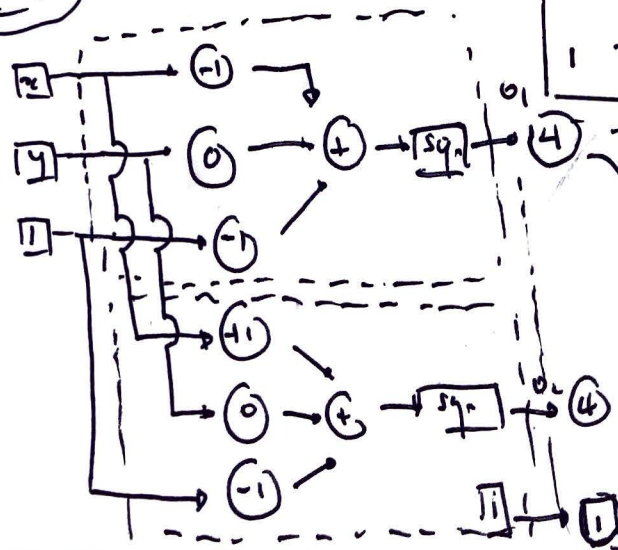
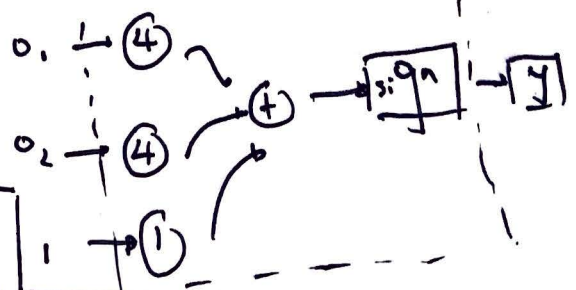
Plane π_2 : $x + z - 2 = 0$
x+z plane:



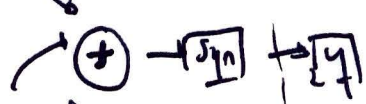
combined:



Goal:



label:



Verification:

$(-2, 0, 1):$

$$o_1 = \text{sqn}(-1 \times \cancel{2}^2 + 0 \times 0 + 1) = \underline{\underline{1}}$$

$$o_2 = \text{sqn}(+1 \times -2 + 0 \times 0 + 1 \times (-1)) = \underline{\underline{-1}}$$

$$y = \text{sqn}(\cancel{4 \times 1} + \cancel{4 \times -1} + 1 \times 1) = \boxed{\underline{\underline{1}}} \quad \checkmark$$

$(2, 0, 1):$

$$o_1 = \text{sqn}(-1 \times 2 + 0 \times 0 - 1) = \underline{\underline{-1}}$$

$$o_2 = \text{sqn}(1 \times \cancel{2}^2 + 0 \times 0 - 1) = \underline{\underline{1}}$$

$$y = \text{sqn}(\cancel{4 \times -1} + \cancel{4 \times 1} + 1) = \boxed{\underline{\underline{1}}} \quad \checkmark$$

$(0, 2, 1):$

$$o_1 = \text{sqn}(-1 \times 0 + 2 \times 0 - 1) = \underline{\underline{-1}}$$

$$o_2 = \text{sqn}(1 \times 0 + 2 \times 0 - 1) = \underline{\underline{-1}}$$

$$y = \text{sqn}(\cancel{4 \times -1} + \cancel{4 \times -1} + 1) = \boxed{\underline{\underline{-1}}} \quad \checkmark$$

$$(0, -2, 1):$$

$$a_1 = \text{sgn}(-1 \times 0 + 1 \times 0 - 1) = \underline{\underline{-1}}$$

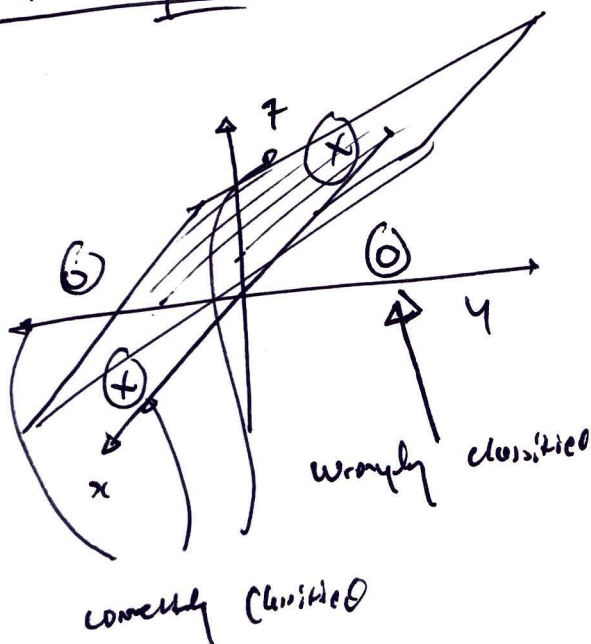
$$a_2 = \text{sgn}(+1 \times 0 + (-2) \times 0 - 1) = \underline{\underline{-1}}$$

$$y = \text{sgn}(-4 - 4 + 1) \left[\begin{matrix} -1 \\ -1 \end{matrix} \right] \checkmark$$

\therefore All correctly classified.

2. If we use one unit, the points cannot be all classified correctly because the cluster points are not linearly separable in \mathbb{R}^3 space with only one decision boundary.
one unit \Rightarrow corresponds to one decision boundary.

Not linearly separable:



\Rightarrow It is impossible to create a ^{single} plane that correctly separates clusters into two distinct categories

\Rightarrow Not linearly separable with one plane/one unit in \mathbb{R}^3 space