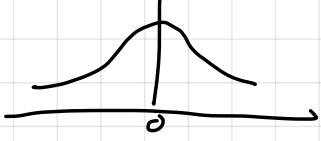


**Problem 1. (15 points).** Consider one layer of a ReLU network. The feature vector is  $d$  dimensional  $\vec{x}$ . The linear transformation is a  $m \times d$  dimensional matrix  $W$ . The output of the ReLU network is a  $m$  dimensional vector  $y$  given by  $\max\{\mathbf{0}, W\vec{x}\}$ . This is a component-wise max function.

- Suppose  $\vec{x}$  is fixed, and all its entries are non-zero.
  - Suppose the entries in  $W$  are all independent, and distributed according to a Gaussian distribution with mean 0, and standard deviation 1 (a  $N(0, 1)$  distribution).
1. Show that the expected number of non-zero entries in the output is  $m/2$ .
  2. Suppose  $\|\vec{x}\|_2^2 = \sigma^2$ , what is the distribution of each entry in  $W\vec{x}$  (the output before applying ReLU function)?
  3. What is the mean of each entry in  $y$  (after ReLU function)?

$$\textcircled{1} \quad \vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}_{d \times 1} \quad W = \begin{bmatrix} w_{11} & \dots & w_{1d} \\ \vdots & \ddots & \vdots \\ w_{m1} & \dots & w_{md} \end{bmatrix}_{m \times d} \quad W\vec{x} = \begin{bmatrix} w_{11} & \dots & w_{1d} \\ \vdots & \ddots & \vdots \\ w_{m1} & \dots & w_{md} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \\ = \begin{bmatrix} x_1 w_{11} + \dots + x_d w_{1d} \\ \vdots \\ x_1 w_{m1} + \dots + x_d w_{md} \end{bmatrix}_{m \times 1}$$

We know that  $W$  has distribution,  
 where mean = 0

$$\text{let } W\vec{x} = q \Rightarrow \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_m \end{bmatrix}_{m \times 1}$$

$$\text{So, } E(q) = 0$$

What is  $E(q_i)$ ? We know that  $\vec{x}$  is fixed or constant.

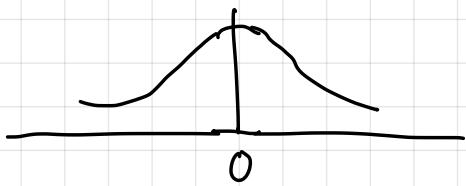
$$\begin{aligned} E(q_i) &= E(x_1 w_{1i} + \dots + x_d w_{di}) \\ &= E(x_1 w_{1i}) + \dots + E(x_d w_{di}) \quad \begin{cases} \text{property of } E(X+Y) = E(X) + E(Y) \\ E(cx) = c E(x) \end{cases} \\ &= x_1 E(w_{1i}) + \dots + x_d E(w_{di}) \\ &= x_1(0) + \dots + x_d(0) \\ &= 0 \end{aligned}$$

→ This is true for all  $E(q_i)$   $i = 1, \dots, m$

$$\text{So, } E(q) = E(W\vec{x}) = 0$$

Since  $W$  has a Gaussian dist.

$W\bar{x}$  where  $\bar{x}$  is constant also has Gaussian dist



mean / expected value = 0

From plot above, we can see that half  $> 0$   
half  $< 0$ .

$$Y = \max \{0, W\bar{x}\} . \quad \text{If } W\bar{x} > 0 , \quad Y > 0 \\ \text{If } W\bar{x} < 0 , \quad Y = 0$$

there's half chance that  $W\bar{x} > 0$  and  $Y > 0$

# of entries

So, the expected value of  $Y > 0$ , is  $\frac{1}{2} \times m^{\leftarrow}$

$$= \boxed{\frac{m}{2}}$$

$$\textcircled{2} \quad \sigma^2 = \|\bar{x}\|_2^2 = x_1^2 + x_2^2 + \dots + x_d^2 \quad \begin{matrix} \text{mean} \\ \uparrow \end{matrix} \quad \begin{matrix} \text{variance} \\ \overbrace{\quad \quad \quad}^d \end{matrix}$$

$$\text{We know } W \sim N(0, 1) = N(\mu, \sigma^2)$$

Based on properties :

$$\mathbb{E}(cx) = c\mathbb{E}(x)$$

so,  $\mathbb{E}(Wx) = x\mathbb{E}(w)$ , since  $x$  is fixed/constant

$$V(cx) = c^2 V(x)$$

$$\text{so, } V(Wx) = x^2 V(w)$$

$$Wx \sim N(\mu, \sigma^2)$$

Find the distribution of each entry of  $W\bar{x} \rightarrow q_1, q_2 \dots, q_m$

$$q_1 = x_1 W_{11} + \dots + x_d W_{1d}$$

$$W_{ii} \bar{x}_i = x_1 W_{11} + \dots + x_d W_{1d}$$

$W_{ii} \bar{x}_i \sim N(\mu, \sigma^2) \leftarrow$  find the distribution

From Q1, we have this proof:

$$\begin{aligned} E(q_1) &= E(x_1 W_{11} + \dots + x_d W_{1d}) \\ &= E(x_1 W_{11}) + \dots + E(x_d W_{1d}) \\ &= x_1 E(W_{11}) + \dots + x_d E(W_{1d}) \\ &= x_1(0) + \dots + x_d(0) \\ &= \underline{\underline{0}} \end{aligned}$$

$\rightarrow$  this is the case for all  $E(q_i)$ ,  $i = 1, \dots, m$

Let's find similar proof for Variance

$$\begin{aligned} V(q_1) &= V(x_1 W_{11} + \dots + x_d W_{1d}) \\ &= V(x_1 W_{11}) + \dots + V(x_d W_{1d}) \quad \leftarrow \text{since entries of } W \text{ are all independent.} \\ &= x_1^2 V(W_{11}) + \dots + x_d^2 V(W_{1d}) \quad \leftarrow V(cx) = c^2 V(x) \text{ and the } x \text{'s in our application is constant.} \\ &= x_1^2(1) + \dots + x_d^2(1) \\ &= x_1^2 + x_2^2 + \dots + x_d^2 \\ &= \underline{\underline{\sigma^2}} \end{aligned}$$

$\uparrow$  this is the case for  $V(q_1), V(q_2), \dots, V(q_m)$

So, distribution for each entry of  $Wx \rightarrow$

$$N(0, \sigma^2)$$

$$\textcircled{3} \quad \hat{y} = \begin{bmatrix} \max(0, x_1 w_{11} + \dots + x_d w_{1d}) \\ \vdots \\ \max(0, x_1 w_{m1} + \dots + x_d w_{md}) \end{bmatrix} = \begin{bmatrix} \max(0, q_1) \\ \vdots \\ \max(0, q_m) \end{bmatrix}$$

Find  $\mu$  of each entry in  $\hat{y}$ :

We have a Gaussian dist for  $y$ , since it is a linear comb of gaussian and linear comb of gaussian is a gaussian by definition.

$$\mu = E(x) = \int_{-\infty}^{\infty} x f(x) dx, \quad f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$E(y_i) = E(\max(0, q_i)) = \int_{-\infty}^{\infty} \max(0, q_i) p(q_i) dq_i$$

$$= \int_{-\infty}^{\infty} \max(0, q_i) \frac{e^{-(q_i - \mu)^2 / 2\sigma^2}}{\sigma \sqrt{2\pi}} dq_i$$

$$= \int_0^{\infty} q_i \frac{e^{-q_i^2 / 2\sigma^2}}{\sigma \sqrt{2\pi}} dq_i$$

$\leftarrow$  bound change to  $(0, \infty)$  because when  $q_i < 0$ , value = 0

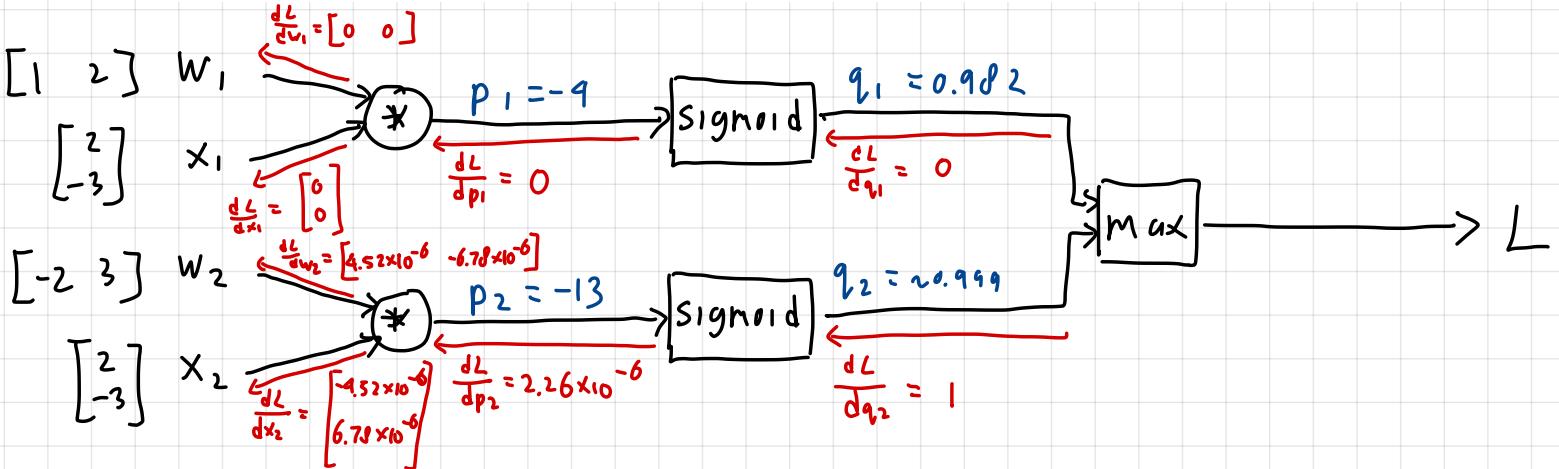
$$= 0 - \left( - \frac{\sigma \Gamma(1)}{\sqrt{2\pi}} \right)$$

$$\mu = \underline{\frac{\sigma}{\sqrt{2\pi}}}$$

**Problem 2. (10 points).** Consider the setting as in the previous problem, with  $m = 2$ , and  $d = 2$ . Let

$$W = \begin{bmatrix} 1 & 2 \\ -2 & 3 \end{bmatrix}, \vec{x} = \begin{bmatrix} 2 \\ -3 \end{bmatrix}. \quad \sigma = \frac{1}{1+e^{-x}}$$

Consider the function  $L = \max\{\sigma(W_{(1)}\vec{x}), \sigma(W_{(2)}\vec{x})\}$ , where  $\sigma$  is the Sigmoid function and  $W_{(i)}$  denotes the  $i$ th row of  $W$ . Please draw the computational graph for this function, and compute the gradients (which will be Jacobians at some nodes!).



Forward :

$$p_1 = [1 \ 2] \cdot \begin{bmatrix} 2 \\ -3 \end{bmatrix} = 2 + (-6) = -4$$

$$p_2 = [-2 \ 3] \cdot \begin{bmatrix} 2 \\ -3 \end{bmatrix} = -4 + (-9) = -13$$

$$q_1 = \frac{1}{1+e^{-4}} = 0.98201379$$

$$q_2 = \frac{1}{1+e^{-13}} = 0.99999773$$

$$L = \max(q_1, q_2) = q_2 = 0.99999773$$

Backward :

$$\frac{dL}{dq_1} = 0, \text{ because } q_2 > q_1$$

$$\frac{dL}{dq_2} = 1, \text{ because } q_2 > q_1$$

$$\begin{aligned} \frac{dL}{dp_1} &= \frac{dL}{dq_1} \cdot \frac{dq_1}{dp_1} = 0 \cdot (\sigma(q_1)(1-\sigma(q_1))) = 0 \cdot \left(\frac{1}{1+e^{-4}} \cdot \left(1 - \frac{1}{1+e^{-4}}\right)\right) \\ &= 0 \cdot 0.01766 = 0 \end{aligned}$$

$$\frac{dL}{dp_2} = \frac{dL}{dq_2} \cdot \frac{dq_2}{dp_2} = 1 \cdot (\sigma(q_2)(1-\sigma(q_2))) = 1 \cdot \left(\frac{1}{1+e^{-13}} \cdot \left(1 - \frac{1}{1+e^{-13}}\right)\right) = 2.26 \times 10^{-6}$$

$$\begin{aligned} \frac{dL}{dw_1} &= \frac{dL}{dp_1} \cdot \frac{dp_1}{dw_1} = \frac{dL}{dp_1} \cdot \left[ \frac{dp_1}{dw_{1(1)}} \quad \frac{dp_1}{dw_{1(2)}} \right] = 0 \cdot \left[ X_{1(1)} \quad X_{1(2)} \right] \\ &= 0 \cdot \begin{bmatrix} 2 & -3 \end{bmatrix} = \begin{bmatrix} 0 & 0 \end{bmatrix} \end{aligned}$$

same dim to  $w_{(1)}$  ↵

$$\frac{dL}{dx_1} = \frac{dL}{dp_1} \cdot \frac{dp_1}{dx_1} = \frac{dL}{dp_1} \cdot \begin{bmatrix} \frac{dp_1}{dx_1(1)} \\ \frac{dp_1}{dx_1(2)} \end{bmatrix} = 0 \cdot \begin{bmatrix} w_{1(1)} \\ w_{1(2)} \end{bmatrix} = 0 \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \boxed{\begin{bmatrix} 0 \\ 0 \end{bmatrix}}$$

↑ Same dim to  $x_1$

$$\frac{dL}{dw_2} = \frac{dL}{dp_2} \cdot \frac{dp_2}{dw_2} = \frac{dL}{dp_2} \cdot \begin{bmatrix} \frac{dp_2}{dw_2(1)} & \frac{dp_2}{dw_2(2)} \end{bmatrix} = 2.26 \times 10^{-6} \cdot \begin{bmatrix} x_{2(1)} & x_{2(2)} \end{bmatrix}$$

$$= 2.26 \times 10^{-6} \cdot \begin{bmatrix} 2 & -3 \end{bmatrix}$$

*Same dim to  $w_{2(2)}$*

$$\leftarrow 1 \times 2 \quad \leftarrow = \boxed{\begin{bmatrix} 4.52 \times 10^{-6} & -6.78 \times 10^{-6} \end{bmatrix}}$$

$$\frac{dL}{dx_2} = \frac{dL}{dp_2} \cdot \frac{dp_2}{dx_2} = \frac{dL}{dp_2} \cdot \begin{bmatrix} \frac{dp_2}{dx_2(1)} \\ \frac{dp_2}{dx_2(2)} \end{bmatrix} = 2.26 \times 10^{-6} \cdot \begin{bmatrix} w_{2(1)} \\ w_{2(2)} \end{bmatrix} = 2.26 \times 10^{-6} \cdot \begin{bmatrix} -2 \\ 3 \end{bmatrix}$$

$$= \boxed{\begin{bmatrix} -4.52 \times 10^{-6} \\ 6.78 \times 10^{-6} \end{bmatrix}}$$

*same dim to  $x_2$*

**Problem 3. (10 points).** Given inputs  $z_1, z_2 \in \mathbb{R}$ , the softmax function is the following:

$$\hat{y} = \frac{e^{z_1}}{e^{z_1} + e^{z_2}}.$$

Let  $y \in \{0, 1\}$ , then define the cross-entropy loss between  $y$  and  $\hat{y}$  be

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

Prove that:

$$\frac{\partial L(y, \hat{y})}{\partial z_1} = \hat{y} - y, \quad \frac{\partial L(y, \hat{y})}{\partial z_2} = y - \hat{y}.$$

$$\text{Proof } \frac{\partial L(y, \hat{y})}{\partial z_1} = \hat{y} - y$$

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

$$\begin{aligned} \frac{\partial L(y, \hat{y})}{\partial z_1} &= \frac{\partial}{\partial z_1} (-y \log(\hat{y})) - \frac{\partial}{\partial z_1} (1-y) \log(1-\hat{y}) \\ &= -y \frac{\partial}{\partial z_1} \left( \log \left( \frac{e^{z_1}}{e^{z_1} + e^{z_2}} \right) \right) - (1-y) \frac{\partial}{\partial z_1} \left( \log \left( 1 - \frac{e^{z_1}}{e^{z_1} + e^{z_2}} \right) \right) \\ &= -y \frac{\partial}{\partial z_1} \left( \log \left( \frac{e^{z_1}}{e^{z_1} + e^{z_2}} \right) \right) - (1-y) \frac{\partial}{\partial z_1} \left( \log \left( \frac{e^{z_2}}{e^{z_1} + e^{z_2}} \right) \right) \\ \frac{\partial}{\partial z_1} \log(e^{z_1} + e^{z_2}) &= -y \frac{\partial}{\partial z_1} \left( \log(e^{z_1}) - \log(e^{z_1} + e^{z_2}) \right) - (1-y) \frac{\partial}{\partial z_1} \left( \log(e^{z_2}) - \log(e^{z_1} + e^{z_2}) \right) \\ = \frac{1}{e^{z_1} + e^{z_2}} (e^{z_1})' &\left( \begin{aligned} &= -y \frac{\partial}{\partial z_1} (z_1 - \log(e^{z_1} + e^{z_2})) - (1-y) \frac{\partial}{\partial z_1} (z_2 - \log(e^{z_1} + e^{z_2})) \\ &= -y \left( \frac{\partial}{\partial z_1} z_1 - \frac{\partial}{\partial z_1} \log(e^{z_1} + e^{z_2}) \right) - (1-y) \left( \frac{\partial}{\partial z_1} z_2 - \frac{\partial}{\partial z_1} \log(e^{z_1} + e^{z_2}) \right) \\ &= -y (1 - \hat{y}) - (1-y) (0 - \hat{y}) \\ &= -y + y\hat{y} + \hat{y} - y\hat{y} \\ &= \hat{y} - y \end{aligned} \right) \quad (\text{proven}) \end{aligned}$$

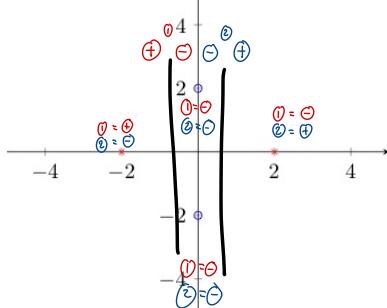
$$P_{\text{root}} \frac{\partial L(y, \hat{y})}{\partial z_2} = y - \hat{y}$$

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

$$\begin{aligned} \frac{\partial L(y, \hat{y})}{\partial z_2} &= \frac{\partial}{\partial z_2} (-y \log(\hat{y})) - \frac{\partial}{\partial z_2} (1-y) \log(1-\hat{y}) \\ &= \frac{\partial}{\partial z_2} (-y \log(\hat{y})) - \frac{\partial}{\partial z_2} (1-y) \log(1-\hat{y}) \\ &= -y \frac{\partial}{\partial z_2} \left( \log \left( \frac{e^{z_1}}{e^{z_1} + e^{z_2}} \right) \right) - (1-y) \frac{\partial}{\partial z_2} \left( \log \left( 1 - \frac{e^{z_1}}{e^{z_1} + e^{z_2}} \right) \right) \\ &= -y \frac{\partial}{\partial z_2} \left( \log \left( \frac{e^{z_1}}{e^{z_1} + e^{z_2}} \right) \right) - (1-y) \frac{\partial}{\partial z_2} \left( \log \left( \frac{e^{z_2}}{e^{z_1} + e^{z_2}} \right) \right) \\ &= -y \frac{\partial}{\partial z_2} (\log(e^{z_1}) - \log(e^{z_1} + e^{z_2})) - (1-y) \frac{\partial}{\partial z_2} (\log(e^{z_2}) - \log(e^{z_1} + e^{z_2})) \\ &= -y \frac{\partial}{\partial z_2} (z_1 - \log(e^{z_1} + e^{z_2})) - (1-y) \frac{\partial}{\partial z_2} (z_2 - \log(e^{z_1} + e^{z_2})) \\ &= -y \left( \frac{\partial}{\partial z_2} z_1 - \frac{\partial}{\partial z_2} \log(e^{z_1} + e^{z_2}) \right) - (1-y) \left( \frac{\partial}{\partial z_2} z_2 - \frac{\partial}{\partial z_2} \log(e^{z_1} + e^{z_2}) \right) \\ &= -y (0 - (1-\hat{y})) - (1-y) (1 - (1-\hat{y})) \\ &= -y(-1 + \hat{y}) - (1-y)(\hat{y}) \\ &= y - y\hat{y} - \hat{y} + y\hat{y} \\ &= y - \hat{y} \quad (\text{proven}) \end{aligned}$$

**Problem 4. (15 points).** Consider datapoints in Figure ??:  $(-2, 0), (2, 0)$  are crosses, and  $(0, 2), (0, -2)$  are circles. Let the crosses be labeled  $+1$ , and the circles be labeled  $-1$ . In this problem the goal

① Separate using 2  $\rightarrow$  boundaries



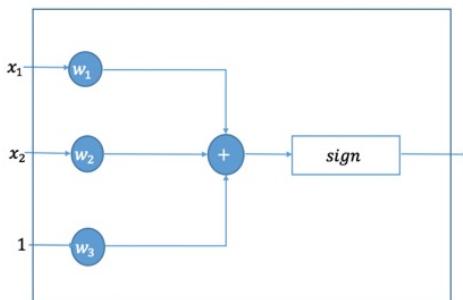
← ② put the result from ①, then classify again using 1 boundary

Figure 1: Neural Networks

is to design a neural network with no error on this dataset.

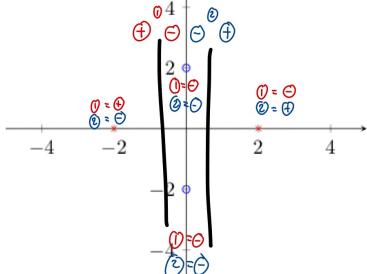
To make things simple, consider the following generalization. We first append a  $+1$  to each input and form a new dataset as follows:  $(-2, 0, 1), (2, 0, 1)$  are labeled  $+1$ , and  $(0, 2, 1), (0, -2, 1)$  are labeled  $-1$ . Note that the last feature is redundant.

We consider the following basic units for our neural networks: Linear transformation followed by hard thresholding. Each unit has three parameters  $w_1, w_2, w_3$ . The output of the unit is the sign of the inner product of the parameters with the input.



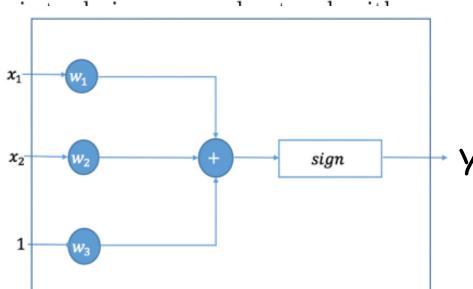
1. Design a neural network with these units that make no error on the datapoints above. (Hint: You can take two units in the first layer, and one in the output layer, a total of three units).
2. Show that if you design a neural network with ONLY one such unit, then the points cannot be all classified correctly.

① Separate using 2  $\rightarrow$  boundaries

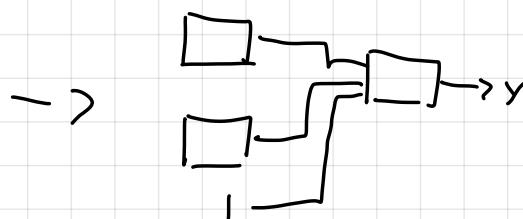


← ② put the result from ①, then classify again using 1 boundary

Figure 1: Neural Networks



$$Y = \text{Sign}(x_1 w_1 + x_2 w_2 + w_3)$$



$$\bar{x} \in \mathbb{R}^3, \{x_1 - x_2, x_3 = 1\}$$

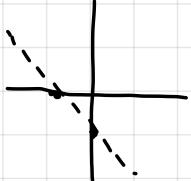
Plane 1 :  $x_1 \geq -1 \cdot x_3$ , where  $x_1 > -1 \Rightarrow \oplus$  and  $x_1 < -1 \Rightarrow \ominus$

$$So, w = \left\{ \frac{1}{w_1}, \frac{0}{w_2}, \frac{-1}{w_3} \right\}$$

Plane 2 :  $x_1 \geq 1 \cdot x_3 = 1$ , where  $x_1 > 1 \Rightarrow \oplus$  and  $x_1 < 1 \Rightarrow \ominus$

$$So, w = \{ -1, 0, -1 \}$$

Plane 3 :



$$x_1 = -1 \Rightarrow x_2 = 0$$

$$\frac{x_1 + 1}{0 + 1} = \frac{x_2 - 0}{-1 - 0}$$

$$x_1 = 0 \Rightarrow x_2 = -1$$

$$\frac{x_1 + 1}{-1 + 1} = \frac{x_2}{-1}$$

$$-1 - x_1 = x_2$$

$$w = \{ 1, 1, 1 \}$$

Test our

$$(0, 2, 1) : y = \underline{\underline{-1}}$$

$$(-2, 0, 1) : y = \underline{\underline{+1}}$$

$$q_1 = \text{sign}(0 \cdot 1 + 2 \cdot 0 + 1 \cdot -1) = -1$$

$$q_1 = \text{sign}(-2 \cdot 1 + 0 \cdot 0 + 1 \cdot -1) = -1$$

$$q_2 = \text{sign}(0 \cdot -1 + 2 \cdot 0 + 1 \cdot -1) = -1$$

$$q_2 = \text{sign}(-2 \cdot -1 + 0 \cdot 0 + 1 \cdot -1) = 1$$

$$Y = \text{Sign}(-1 \cdot 1 + -1 \cdot 1 + 1 \cdot 1) = -1$$

$$Y = \text{Sign}(-1 \cdot 1 + 1 \cdot 1 + 1 \cdot 1) = 1$$

$$(0, -2, 1) : y = \underline{\underline{-1}}$$

$$(2, 0, 1) : y = \underline{\underline{+1}}$$

$$q_1 = \text{sign}(0 \cdot 1 + -2 \cdot 0 + 1 \cdot -1) = -1$$

$$q_1 = \text{sign}(2 \cdot 1 + 0 \cdot 0 + 1 \cdot -1) = 1$$

$$q_2 = \text{sign}(0 \cdot -1 + -2 \cdot 0 + 1 \cdot -1) = -1$$

$$q_2 = \text{sign}(2 \cdot -1 + 0 \cdot 0 + 1 \cdot -1) = -1$$

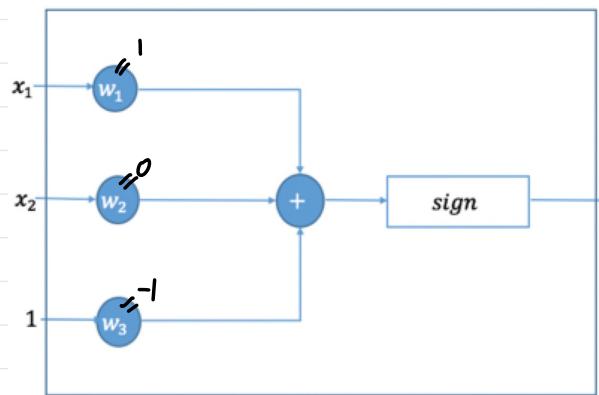
$$Y = \text{Sign}(-1 \cdot 1 + -1 \cdot 1 + 1 \cdot 1) = -1$$

$$Y = \text{Sign}(1 \cdot 1 + -1 \cdot 1 + 1 \cdot 1) = 1$$

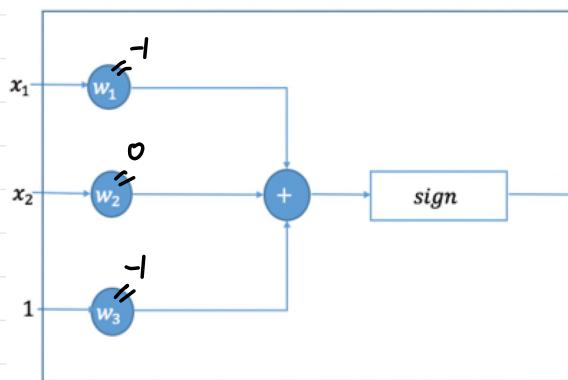
v

v

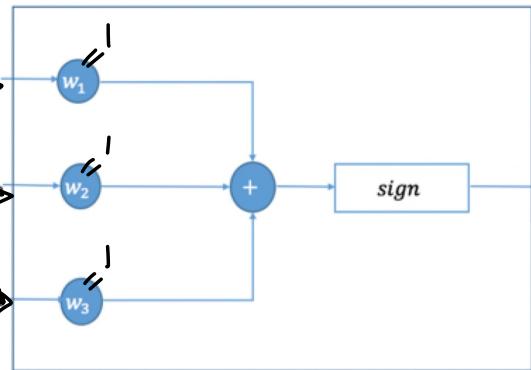
No errors!



$q_1$

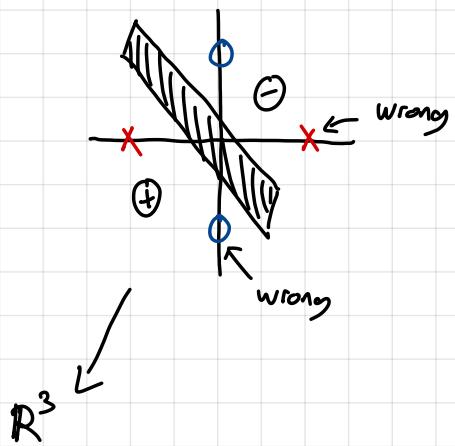


$q_2$



label/  
 $\hat{y}$

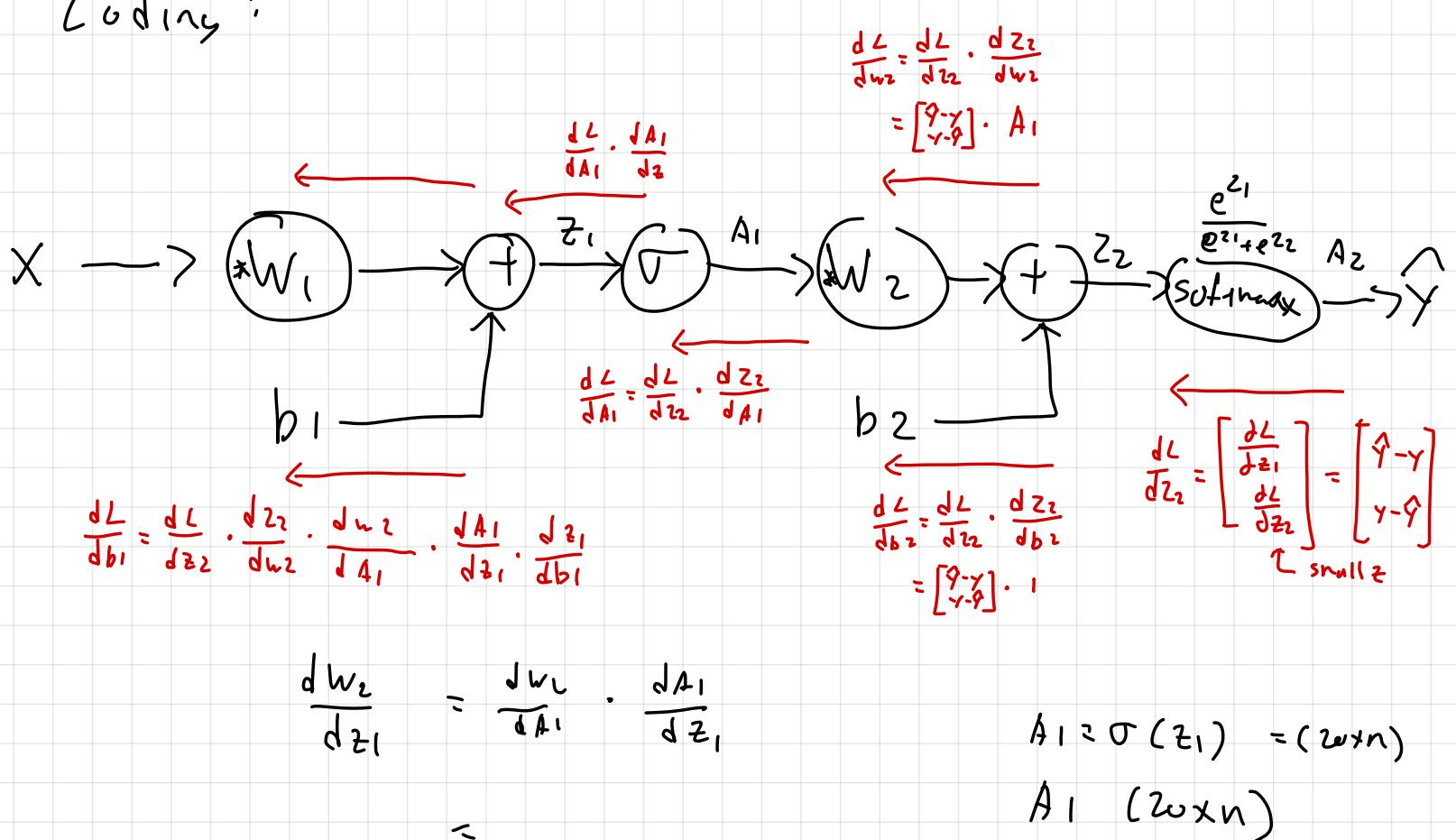
②



→ there is no single 3D hyperplane that can correctly classify the 4 points, as the points cannot be classified with a single boundary. As shown in  $Q_1$  explanation, we need at least 2 boundaries, which can be achieved using 2 boundaries.

Since using a single unit means a single boundary, so the points cannot be correctly classified. Thus, we need more than one unit.

Coding:



$$\frac{dL}{dw_1} = \frac{dL}{dz_2} \cdot \frac{dA_1}{dz_2} \cdot \frac{dt_1}{dz_1}$$

$2w \times n \quad \cdot \quad 2w \times n$

