

Q1

$$d = 0.45 \sqrt{h}$$

$$y = w_0 v + w_1 h$$

$$w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$x = \begin{bmatrix} v \\ h \end{bmatrix}$$

$$y_i = w^T x_i + b$$

I. why is $\hat{R}(w) = \frac{1}{I} \|X_w - d\|^2$?

$$* R(w) = E(l(n, w)) = \frac{1}{I} \sum_i l(n_i, w) = \frac{1}{I} \sum_i l(f_{n_i}(w), y_i)$$

$$L(y, d) = (y - d)^2$$

$$\Leftrightarrow L(y, d) = (w_0 v + w_1 h - d)^2$$

$$R(w) = \frac{1}{I} \sum_i (w_0 v_i + w_1 h_i - d_i)^2 = \frac{1}{I} \sum_i (w^T n_i - d_i)^2$$

$$= \frac{1}{I} \left((w^T n_1 - d_1)^2 + (w^T n_2 - d_2)^2 + \dots + (w^T n_I - d_I)^2 \right)$$

$$= \frac{1}{I} \left([(w^T n_1 - d_1) \dots (w^T n_I - d_I)] \underbrace{\begin{bmatrix} w^T n_1 - d_1 \\ \vdots \\ w^T n_I - d_I \end{bmatrix}}_{H} \right)$$

$$= \frac{1}{I} \|H\|^2 = \frac{1}{I} \left\| \begin{bmatrix} w^T n_1 - d_1 \\ \vdots \\ w^T n_I - d_I \end{bmatrix} \right\|_2^2 \text{ H}$$

next page

$$\begin{bmatrix} w^T n_1 - d_1 \\ \vdots \\ w^T n_I - d_I \end{bmatrix} = \begin{bmatrix} w^T n_1 \\ \vdots \\ w^T n_I \end{bmatrix} - \begin{bmatrix} d_1 \\ \vdots \\ d_I \end{bmatrix} = \begin{bmatrix} n_1^T w \\ \vdots \\ n_I^T w \end{bmatrix} - \underbrace{\begin{bmatrix} d_1 \\ \vdots \\ d_I \end{bmatrix}}_d$$

$$H = \begin{bmatrix} n_1^T \\ \vdots \\ n_I^T \end{bmatrix} w - d = Xw - d \quad \text{(b)}$$

$$(a), (b) \Rightarrow R(w) = \frac{1}{I} \|Xw - d\|^2 \quad \text{(c)}$$

II. gradient of risk:

$$\text{gradient } f = \begin{bmatrix} f_1(w) \\ \vdots \\ f_N(w) \end{bmatrix} = \nabla f(w)$$

$$\nabla R(w) = \begin{bmatrix} R'_1(w) \\ \vdots \\ R'_I(w) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial w_0} R(w) \\ \vdots \\ \frac{\partial}{\partial w_I} R(w) \end{bmatrix}$$

$$R(w) = \frac{1}{I} \left((w^T w_0 - d_1)^2 + \dots + (w^T w_0 - d_I)^2 \right)$$

$$R(w) = \frac{1}{I} \sum_i (w^T w_0 - d_i)^2$$

$$\frac{\partial}{\partial w_n} \left(\frac{1}{I} \sum_i (w^T w_0 - d_i)^2 \right) = \frac{1}{I} \sum_i \frac{\partial}{\partial w_n} (w^T w_0 - d_i)^2$$

$$\frac{\partial}{\partial w_n} (w^T w_0 - d_i)^2 = \frac{\partial}{\partial w_n} (w_0 v_i + w_1 h_i - d_i)^2$$

$$\frac{\partial}{\partial w_0} (w_0 v_i + w_1 h_i - d_i)^2 = 2(w_0 v_i + w_1 h_i - d_i)(v_i)$$

$$\frac{\partial}{\partial w_1} (w_0 v_i + w_1 h_i - d_i)^2 = 2(w_0 v_i + w_1 h_i - d_i)(h_i)$$

$\approx \nabla R(w) = \begin{bmatrix} \frac{1}{I} \sum_i 2 v_i (w_0 v_i + w_1 h_i - d_i) \\ \frac{1}{I} \sum_i 2 h_i (w_0 v_i + w_1 h_i - d_i) \end{bmatrix}$

$$w^{(t)} \leftarrow w^{(t-1)} - \eta \begin{bmatrix} \frac{1}{I} \sum_i 2 v_i (w_0^{(t-1)} v_i^{(t-1)} + w_1^{(t-1)} h_i^{(t-1)} - d_i) \\ \frac{1}{I} \sum_i 2 h_i (w_0^{(t-1)} v_i^{(t-1)} + w_1^{(t-1)} h_i^{(t-1)} - d_i) \end{bmatrix}$$

→ update rule

III. find minimizer

$$\nabla R(w) = \begin{bmatrix} \frac{1}{I} \sum_i 2v_i^o (w_0 v_i^o + w_1 h_i^o - d_i) \\ \frac{1}{I} \sum_i 2h_i^o (w_0 v_i^o + w_1 h_i^o - d_i) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\frac{1}{I} \sum_i 2v_i^o (w_0 v_i^o + w_1 h_i^o - d_i) = 0 \quad \sum_i v_i^o (w_0 v_i^o + w_1 h_i^o - d_i) = 0$$

$$\frac{1}{I} \sum_i 2h_i^o (w_0 v_i^o + w_1 h_i^o - d_i) = 0 \quad \rightarrow \sum_i h_i^o (w_0 v_i^o + w_1 h_i^o - d_i) = 0$$

$$\sum_i v_i^{o^2} w_0 + v_i h_i w_1 - v_i d_i = 0$$

(w_0, w_1 unrelated to i)

$$\sum_i v_i h_i w_0 + h_i^{o^2} w_1 - h_i d_i = 0$$

$$(\sum_i v_i^2) w_0 + (\sum_i v_i h_i) w_1 - (\sum_i v_i d_i) = 0$$

$$(\sum_i v_i h_i) w_0 + (\sum_i h_i^{o^2}) w_1 - (\sum_i h_i d_i) = 0$$

$$V w_0 + A w_1 = B$$

$$A w_0 + H w_1 = C$$

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \begin{bmatrix} V & A \\ A & H \end{bmatrix}^{-1} \begin{bmatrix} B \\ C \end{bmatrix}$$

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \frac{1}{VH-A^2} \begin{bmatrix} H & -A \\ -A & V \end{bmatrix} \begin{bmatrix} B \\ C \end{bmatrix}$$

$$M \cdot \text{adj}(M) = \det(M) I$$

$$M^{-1} = \frac{1}{\det(M)} \text{adj}(M)$$

Cofactor(M)^T

$$\begin{bmatrix} +|H| & -|A| \\ -|A| & +|V| \end{bmatrix}^T$$

$$\begin{bmatrix} V & A \\ A & H \end{bmatrix}^{-1} = \frac{1}{VH-A^2} \begin{bmatrix} H & -A \\ -A & V \end{bmatrix}$$

$$\begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = \frac{1}{VH-A^2} \begin{bmatrix} HB - AC \\ -AB + VC \end{bmatrix}$$

$$w_0 = \frac{(\sum h_i^{o^2})(\sum v_i d_i) - (\sum v_i h_i)(\sum h_i d_i)}{(\sum v_i^{o^2})(\sum h_i^{o^2}) - (\sum v_i h_i)^2}$$

$$w_1 = \frac{(\sum v_i^{o^2})(\sum h_i d_i) - (\sum v_i h_i)(\sum v_i d_i)}{(\sum v_i^{o^2})(\sum h_i^{o^2}) - (\sum v_i h_i)^2}$$

IV.

$$w_0 = \frac{(\varepsilon_{h_i^2})(\varepsilon_{v_i h_i}) - (\varepsilon_{v_i h_i})(\varepsilon_{h_i^2})}{(\varepsilon_{v_i^2})(\varepsilon_{h_i^2}) - (\varepsilon_{v_i h_i})^2} \quad \left\{ \text{from III} \right.$$

$$\textcircled{1} w_1 = \frac{(\varepsilon_{v_i^2})(\varepsilon_{h_i^2}) - (\varepsilon_{v_i h_i})(\varepsilon_{v_i h_i})}{(\varepsilon_{v_i^2})(\varepsilon_{h_i^2}) - (\varepsilon_{v_i h_i})^2}$$

$$\begin{aligned} & (v_o^2 + v_n^2)(h_o^2 + h_n^2) - (v_o h_o + v_n h_n)^2 \\ & \cancel{v_o^2 h_o^2} + \cancel{v_o^2 h_n^2} + \cancel{v_n^2 h_o^2} + \cancel{v_n^2 h_n^2} - v_o^2 h_o^2 - v_n^2 h_n^2 - 2 v_o v_n h_o h_n \\ & \boxed{(v_o h_n - v_n h_o)^2} \quad \text{I} \end{aligned}$$

$$\textcircled{2} w_1 = \frac{(v_o^2 + v_n^2)(h_o d_o + h_n d_n) - (v_o h_o + v_n h_n)(v_o d_o + v_n d_n)}{v_o^2 h_o d_o + v_o^2 h_n d_n + v_n^2 h_o d_o + v_n^2 h_n d_n - v_o^2 h_o d_o - v_o v_n h_o d_n - v_o v_n h_n d_o - v_n^2 h_n d_j}$$

$$\textcircled{3} w_1 = \frac{v_o^2 h_p d_n + v_n^2 h_o d_o - v_o v_n (h_o d_n + h_n d_o)}{(v_o h_n - v_n h_o)^2} = \frac{(v_o h_n - v_n h_o)(v_o d_n - v_n d_o)}{(v_o h_n - v_n h_o)^2}$$

$$w_0 = \frac{h_o^2 v_n d_j + h_n^2 v_o d_o - h_o h_n (v_o d_n + v_n d_o)}{(v_o h_n - v_n h_o)^2} = \frac{(v_o h_i - v_n h_o)(h_i d_n - h_o d_o)}{(v_o h_n - v_n h_o)^2}$$

$$\Rightarrow (v_o h_n - v_n h_o) \neq 0 \Rightarrow w_1 = \frac{h_i d_o - h_o d_n}{v_o h_n - v_n h_o} \quad \text{I}$$

* hypothesis: $\det \begin{bmatrix} v_o & h_o \\ v_n & h_n \end{bmatrix} \neq 0$

$$w_0 v_o + w_1 h_o - d_o = \frac{-v_o h_i d_n + v_n h_o d_o + v_o h_o d_n - v_n h_i d_o}{v_o h_n - v_n h_o} - d_o = 0$$

$$w_0 v_n + w_1 h_n - d_j = \frac{-v_n h_i d_n + v_n h_o d_o + v_o h_n d_n - v_n h_i d_o}{v_o h_n - v_n h_o} - d_j = 0$$

Continued on
next page

IV

for $J \geq 3$

$$w_0 = \frac{(\sum h_i^2)(\sum v_i d_i) - (\sum v_i h_i)(\sum h_i d_i)}{(\sum v_i^2)(\sum h_i^2) - (\sum v_i h_i)^2}$$

$$w_n = \frac{(\sum v_i^2)(\sum h_i d_i) - (\sum v_i h_i)(\sum v_i d_i)}{(\sum v_i^2)(\sum h_i^2) - (\sum v_i h_i)^2}$$

$$w_i = \frac{v_i^2 h_1^2 + v_1^2 h_2^2 + v_1^2 h_3^2 - \dots - (v_1 h_1 + v_2 h_2 + v_3 h_3)^2}{v_1^2 h_2^2 + v_1^2 h_3^2 + v_2^2 h_1^2 + v_2^2 h_3^2 + v_3^2 h_1^2 + v_3^2 h_2^2 - 2(v_1 h_1 v_2 h_2 + v_1 h_1 v_3 h_3 + v_2 h_2 v_3 h_3)}$$

$$w_0 = \frac{(v_1 h_2 - v_2 h_1)(h_2 d_1 - h_1 d_2) + (v_1 h_3 - v_3 h_1)(h_3 d_1 - h_1 d_3) + (v_2 h_3 - v_3 h_2)(h_3 d_2 - h_2 d_3)}{(v_1 h_2 - v_2 h_1)^2 + (v_1 h_3 - v_3 h_1)^2 + (v_2 h_3 - v_3 h_2)^2}$$

$$w_1 = \frac{(v_1 h_2 - v_2 h_1)(v_2 d_1 - v_1 d_2) + (v_1 h_3 - v_3 h_1)(v_3 d_1 - v_1 d_3) + (v_2 h_3 - v_3 h_2)(v_3 d_2 - v_2 d_3)}{(v_1 h_2 - v_2 h_1)^2 + (v_2 h_3 - v_3 h_2)^2 + (v_3 h_1 - v_1 h_3)^2}$$

Now we could see that as the number of subpoints is increased, w_1 and w_2 would become more complicated and $\hat{R}_{=0}$ would be almost infeasible and -

- also \hat{R} would increase because we are approximating a non-linear function with a linear function!!

V. for w_t to converge in one step, we -

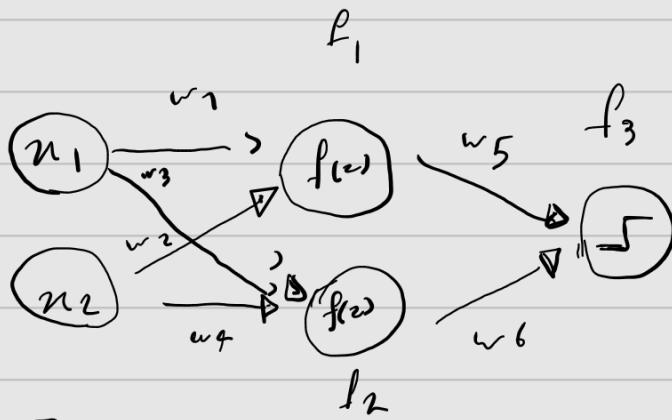
- need $w^{t+1} \leftarrow w^t - \eta \nabla R$ to get to the optimal w in one step, meaning that -
 - both w_0 and w_1 would need to change -
 - to w_0' and w_1' in one step -
- we could notice that because η is applied to both w_0, w_1 then the initial values and their distance from optimal w_0 and w_1 would need to be exactly the same so that a single η -
- would make both optimal in one step -
 - which we know almost never happens!!

VI. If we use $\begin{bmatrix} n_1 \\ 0 \end{bmatrix}$, then both n_1, n_2

could be fine tuned regardless of the other one so that w_0 and w_1 could converge to optimal values in one step. Q.E.D.

Q 2

I.



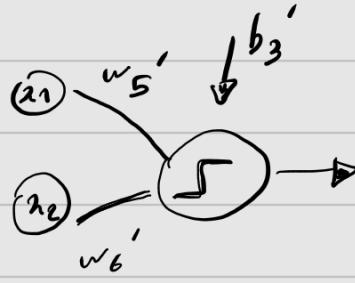
$$f_{(2)} = 2$$

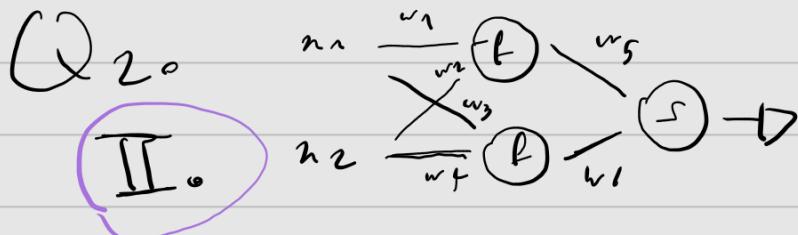
$$f_1 \rightarrow w_1 n_1 + w_2 n_2 + b_1 \quad \rightsquigarrow S((w_1 n_1 + w_2 n_2 + b_1) w_5 + (w_3 n_1 + w_4 n_2 + b_2) w_6 + b_3)$$

$$f_2 \rightarrow w_3 n_1 + w_4 n_2 + b_2$$

$$S((w_1 w_5 + w_3 w_6) n_1 + (w_2 w_5 + w_4 w_6) n_2 + w_5 b_1 + w_6 b_2 + b_3)$$

$$S(w'_5 n_1 + w'_6 n_2 + b'_3) =$$





Let's try to create AND and OR with ReLU

$$\text{AND} \rightarrow w_1=1, w_2=2, b_1=-2$$

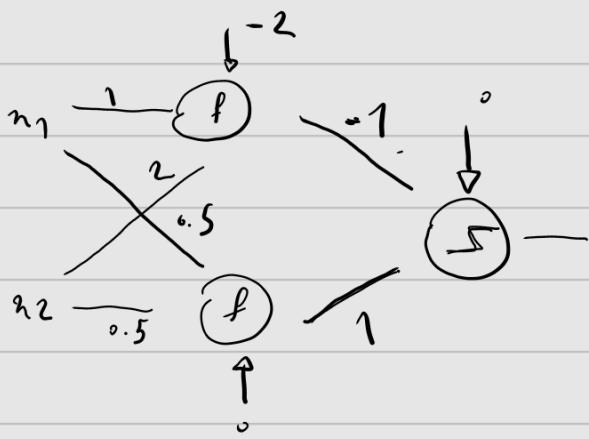
OR ~~ReLU can't~~ but we can create a semi-OR

$$w_3=0.5, w_4=0.5$$

$$b_2=0$$

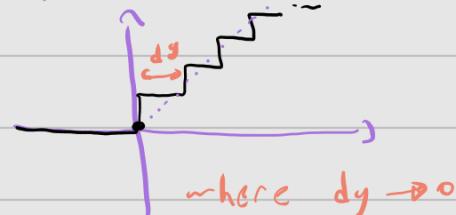
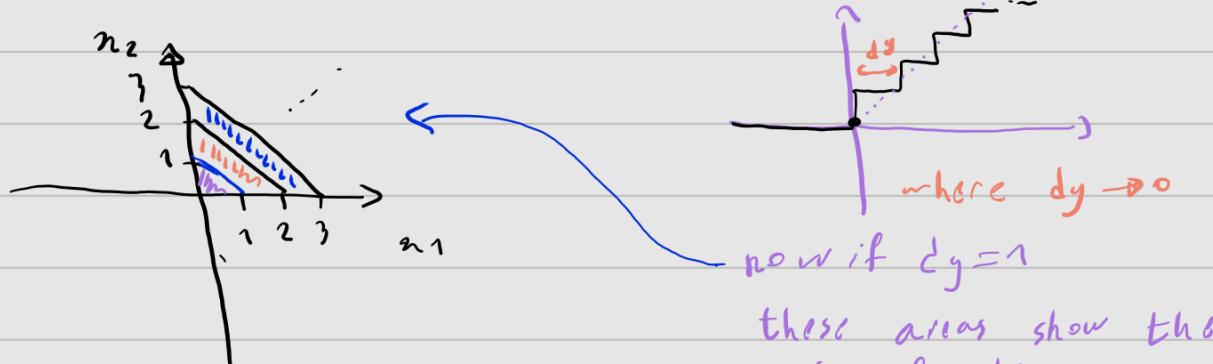
* what we try to

do is we use ReLU as a small classifier like perceptron



* note → geometrically analyzing

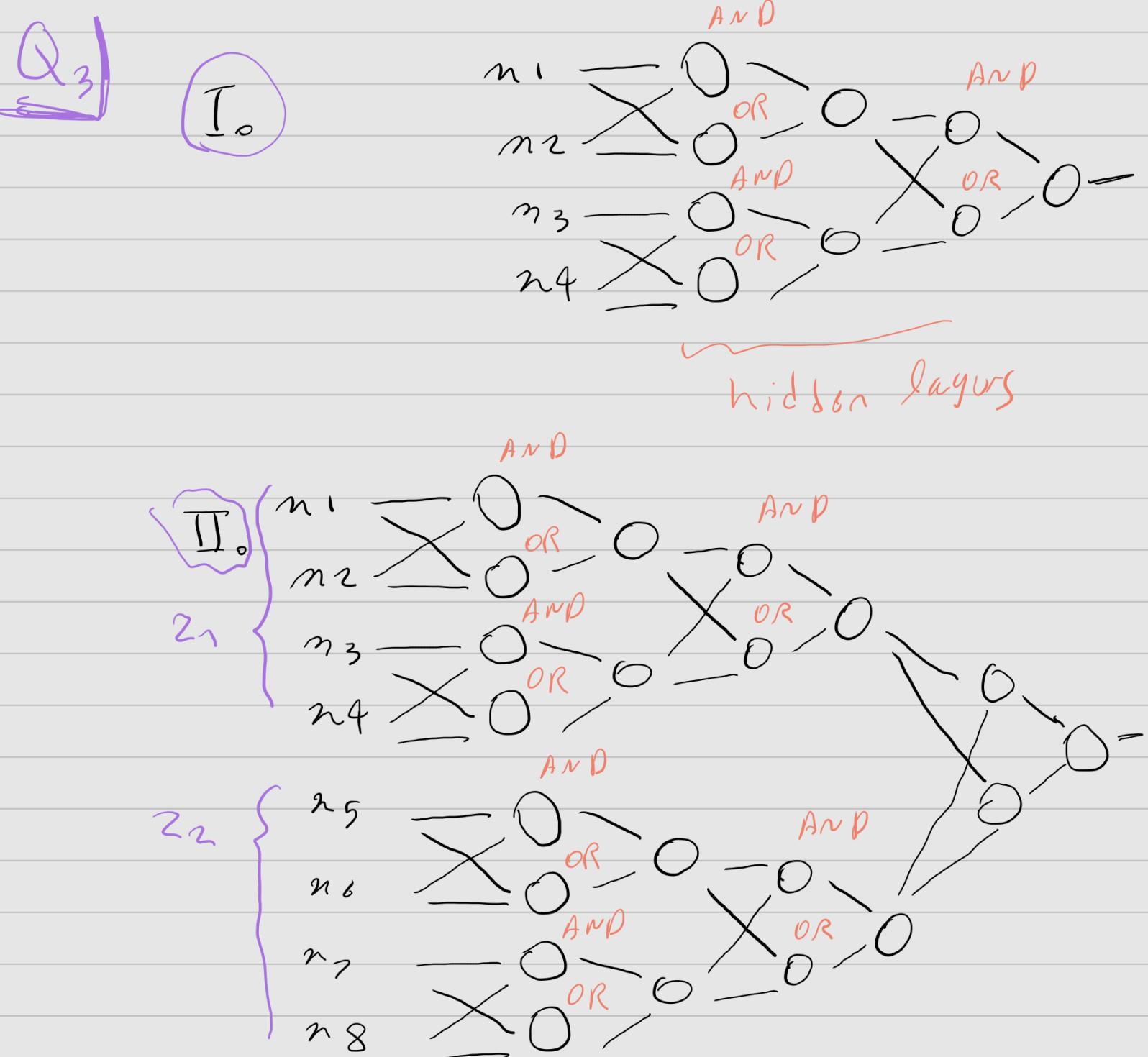
ReLU is basically like this



where $dy \rightarrow 0$

now if $dy = 1$

these areas show the classifications



$$(n_1 \oplus n_2 \oplus n_3 \oplus n_4) \oplus (n_5 \oplus n_6 \oplus n_7 \oplus n_8)$$

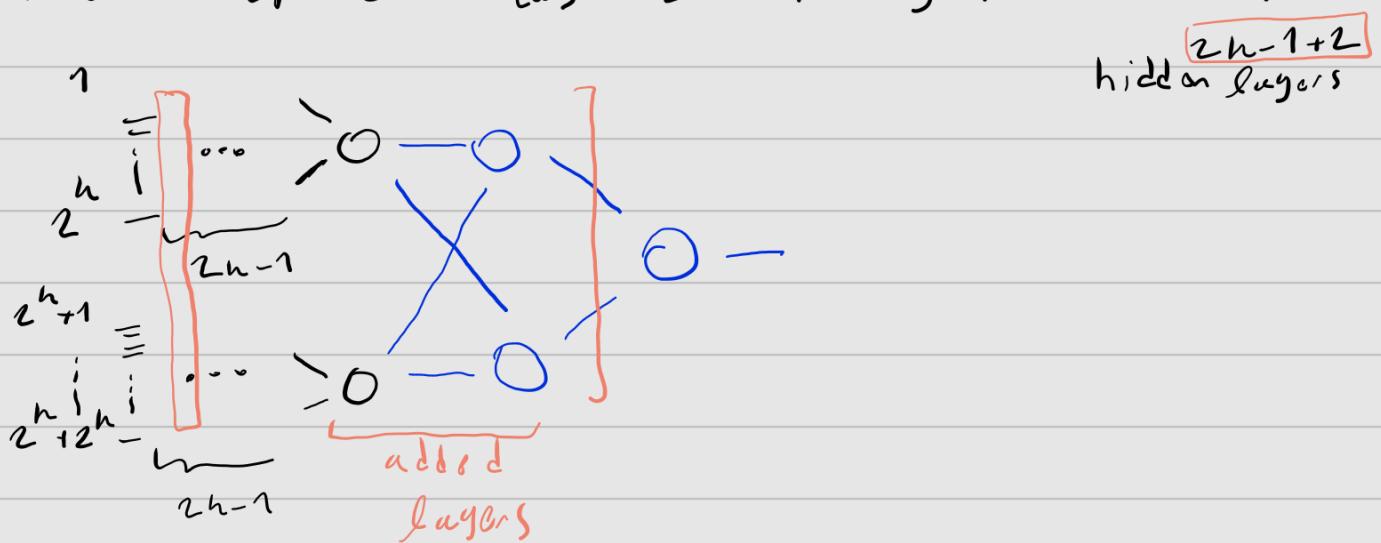
$\underbrace{\qquad\qquad\qquad}_{z_1}$ $\underbrace{\qquad\qquad\qquad}_{z_2}$

III, IV. If input = $N = 2^k$

hidden layers = 2^{k-1}

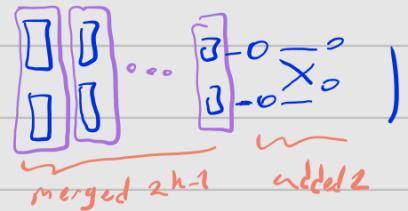
proof: we know the base of induction ($k=1$) is true (1 hidden layer)

induction step: 2^n has 2^{n-1} layers so 2^{n+1} has



$$(n_1 \oplus n_2 \oplus \dots \oplus n_{2^n}) \otimes (n_{2^{n-1}} \oplus \dots \oplus n_{2^{n+1}})$$

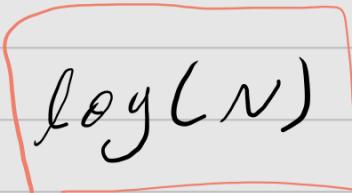
needs 2^{n-1} needs 2^{n-1}
2ⁿ⁻¹ first hidden layers merge



so $n+1$ layers should have

$$2(n+1)-1 = (2n-1) + 2$$

✓ Proof complete

 we could conclude that
having deep-NNs would decrease
the number of needed neurons($=N$)
to a  $\log(N)$ order 