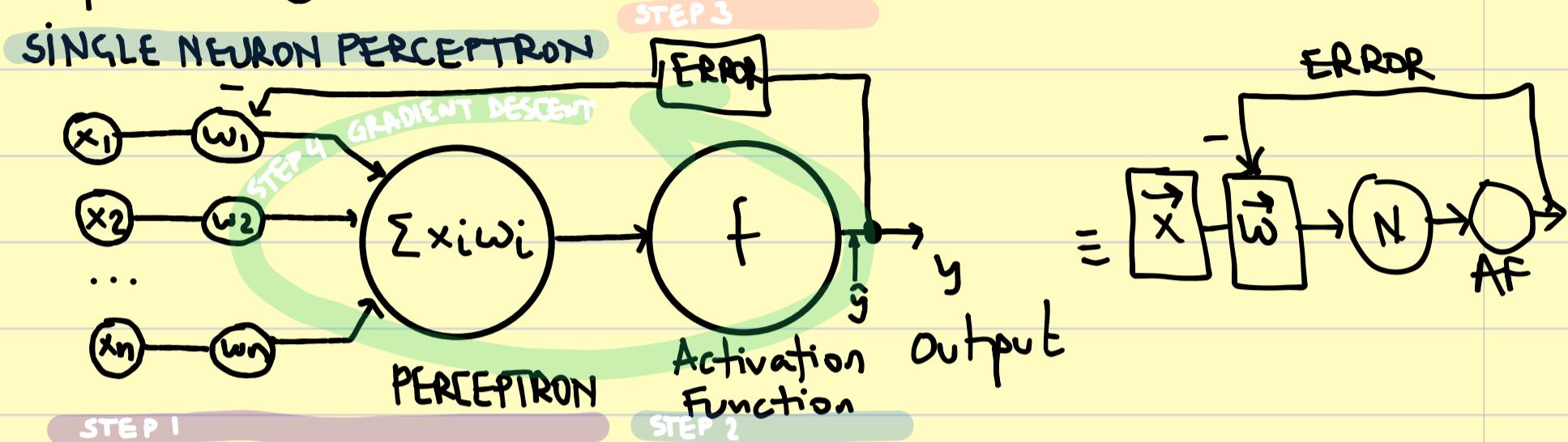


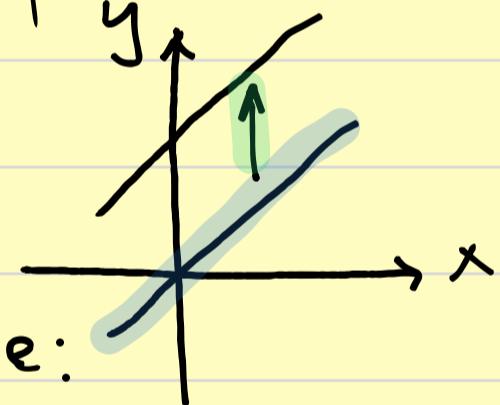
## Deep Learning .. by hand



**Step 1.** Multiply all input values  $x_i$  with their corresponding weights  $w_i$  and then calculate the weighted sum.

**Step 2.** Apply an activation function is applied to the above giving an output of the form:

$$\hat{y} = f\left(\sum x_i w_i + \text{bias}\right)$$



Examples of activation functions are:

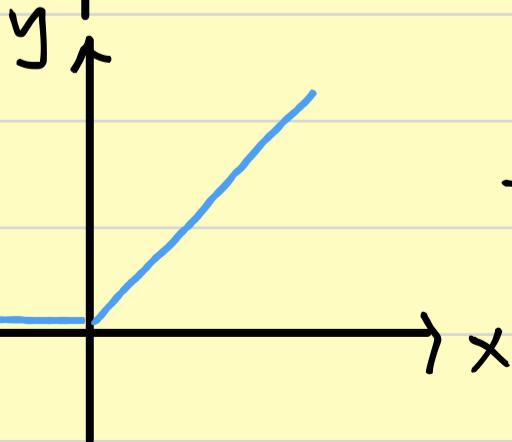
- Step function



$$f(x) = \begin{cases} 0 & x < 0 \\ 1 & x > 0 \end{cases}$$

- Rectified Linear

- Unit (ReLU)



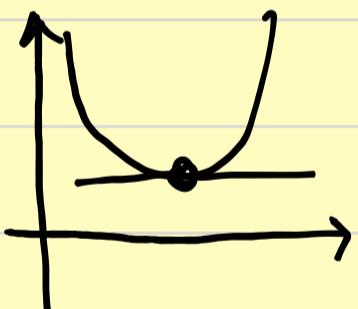
$$f(x) = \begin{cases} 0 & x < 0 \\ x & x > 0 \end{cases}$$

**Step 3.** We need to calculate the weight values which make the equation  $\hat{y} = f(\sum w_i x_i + b)$  as similar as possible to a prediction.

For this we use the quadratic of the difference between the two values  $(y - \hat{y})^2$  = ERROR. This quadratic difference is called the quadratic error:

$$\text{COST FUNCTION} = \frac{1}{2} (y - \hat{y})^2$$

Explanation: the error ( $y - \hat{y}$ ) is squared so that all values are positive (they do not cancel each other).



we divide it by two because we want to find the minimum of the function  $C$  and for this we will find the first derivative of the cost function.

$$C = \frac{1}{2} \left( y - \hat{y} \right)^2 = \frac{1}{2} \left( y - \sum w_i x_i - b \right)^2$$

We call this function error or cost and we want to find the values of  $\vec{w} [w_1, w_2, \dots, w_n]$  that minimize the cost.

$$\frac{dC}{dw_1} = \frac{1}{2} \cdot 2 \cdot [y - w_1x_1 - w_2x_2 - \dots - w_nx_n] \cdot (-x_1)$$

$$\frac{dC}{dw_2} = \frac{1}{2} \cdot 2 \cdot [y - w_1x_1 - w_2x_2 - \dots - w_nx_n] \cdot (-x_2)$$

...

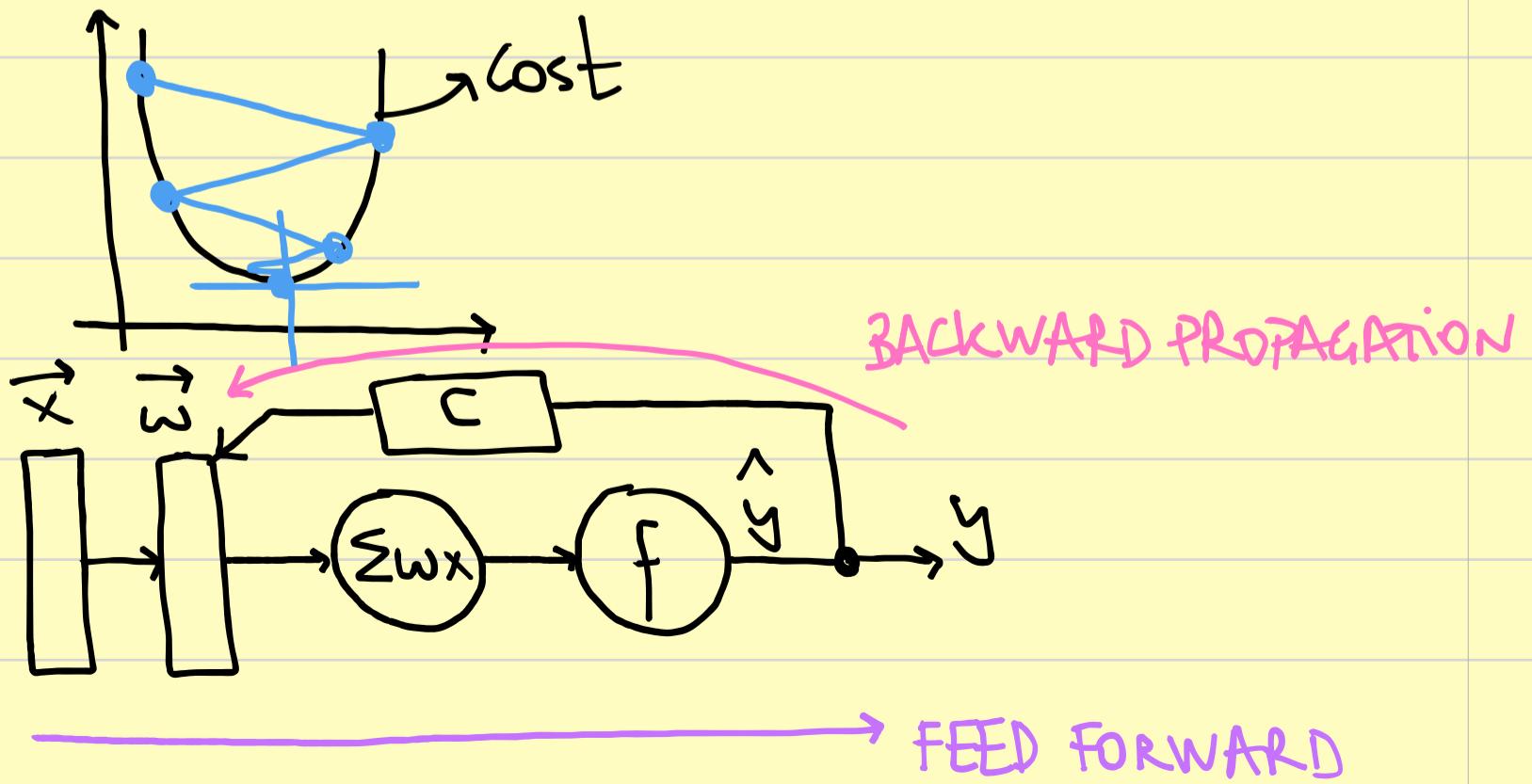
$$\frac{dC}{dw_n} = \frac{1}{2} \cdot 2 \cdot [y - w_1x_1 - w_2x_2 - \dots - w_nx_n] \cdot (-x_n)$$

These derivatives are called gradients.

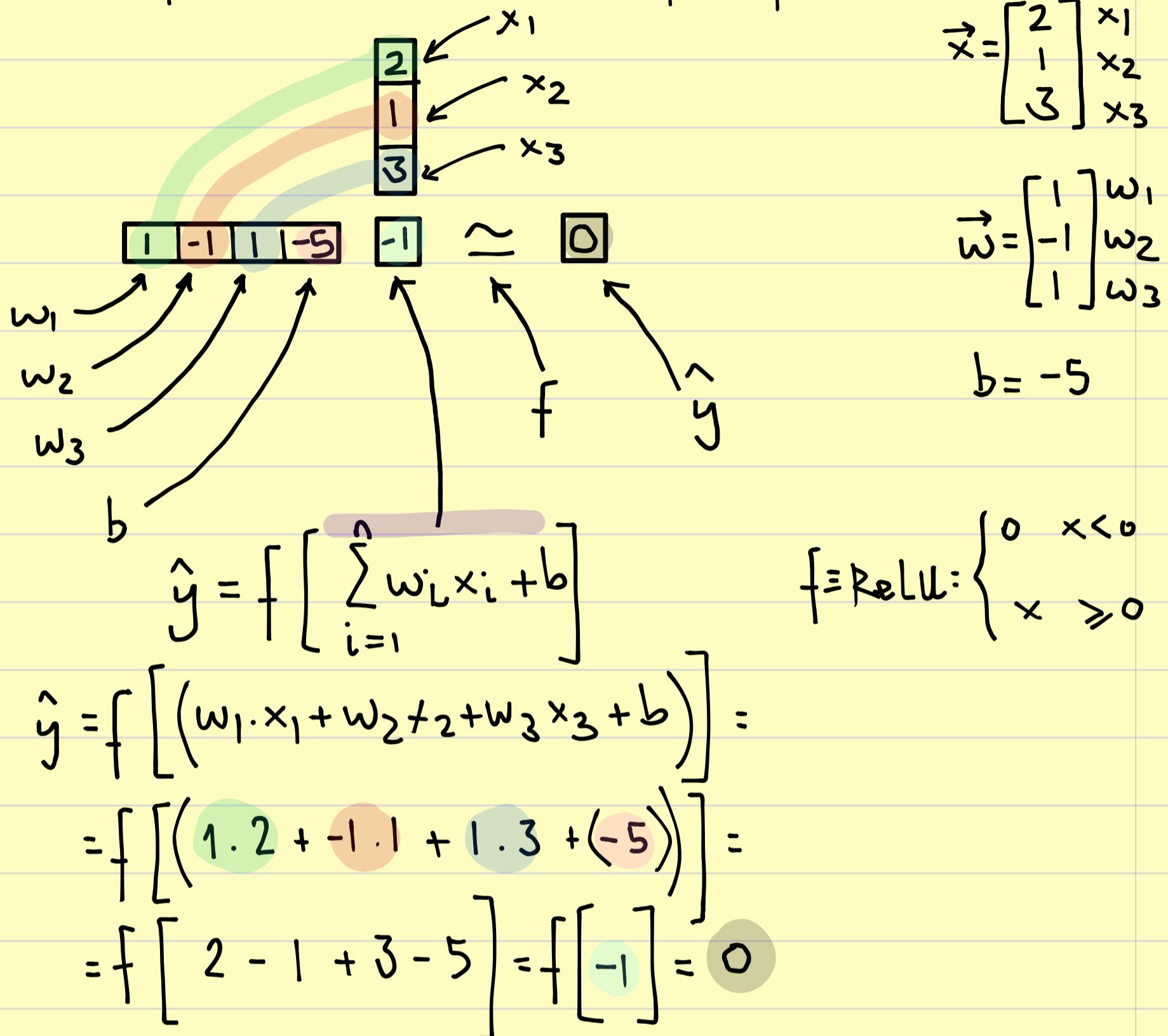
#### Step 4. Gradient descent

In order to reach a minima in the cost function we will start taking small steps in the direction of the minima (opposite direction of the gradient).

We do so by updating the weights by a small amount. This amount is called ..learning rate..



## Example 1. Single neuron perception forward Pass



Example 2. Calculate the output of two layers of FORWARD PASS

perceptrons with 3 and 4 neurons respectively.

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} [1 \times 3]$$

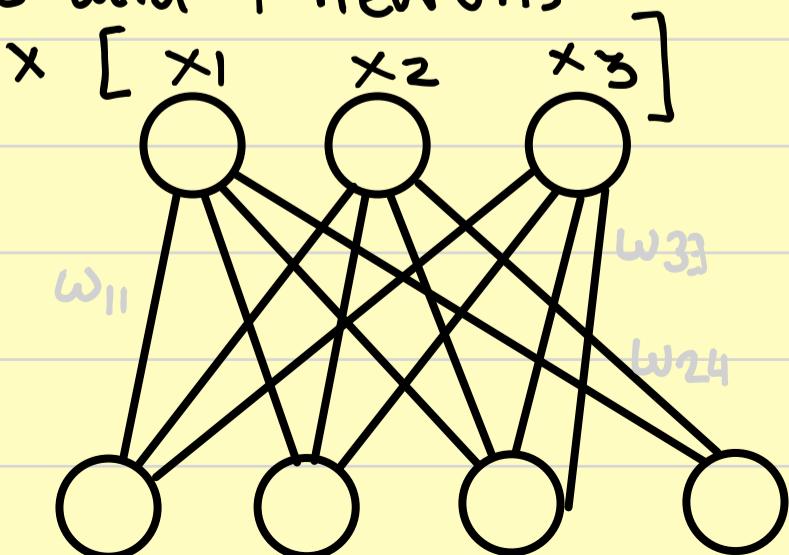
1	-1	1	-5
1	1	0	0
0	1	1	1
1	0	1	-2

$$\vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$\vec{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} 0 \\ 3 \\ 5 \\ 3 \end{bmatrix}$$

$$\approx$$



$$w = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix}$$

$$b = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}$$

$f_{\text{ReLU}}$

$$\hat{y} [ \hat{y}_1 \quad \hat{y}_2 \quad \hat{y}_3 \quad \hat{y}_4 ]$$

Fully connected multi-layer perception.

$[3 \times 4]$

$[1 \times 4]$

$$f(\sum w_i x_i + b) = f\left(\begin{bmatrix} 2 & 1 & 3 \end{bmatrix} \begin{bmatrix} 1 & -1 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \end{bmatrix} + \begin{bmatrix} -5 \\ 0 \\ 1 \\ -2 \end{bmatrix}\right) = f\left(\begin{bmatrix} -1 \\ 3 \\ 5 \\ 3 \end{bmatrix}\right) =$$

$$= \begin{bmatrix} 0 \\ 3 \\ 5 \\ 3 \end{bmatrix} = \hat{y}$$

$2 \cdot 1 + 1 \cdot (-1) + 3 \cdot 1 - 5 = -1$   
 $2 \cdot 1 + 1 \cdot 1 + 3 \cdot 0 + 0 = 3$   
 $2 \cdot 0 + 1 \cdot 1 + 3 \cdot 1 + 1 = 5$   
 $2 \cdot 1 + 1 \cdot 0 + 3 \cdot 1 - 2 = 3$

**Example 3.** 3 layer perceptron (forward pass)  $[3, 4, 2]$

HIDDEN LAYER

$$[3 \times 4] \quad [1 \times 4]$$

$$w_{1h} \quad b_1$$

$$\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \rightarrow \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix} \hat{h}$$

$$\begin{bmatrix} 1 & -1 & 1 & -5 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & -2 \end{bmatrix} [4 \times 2]$$

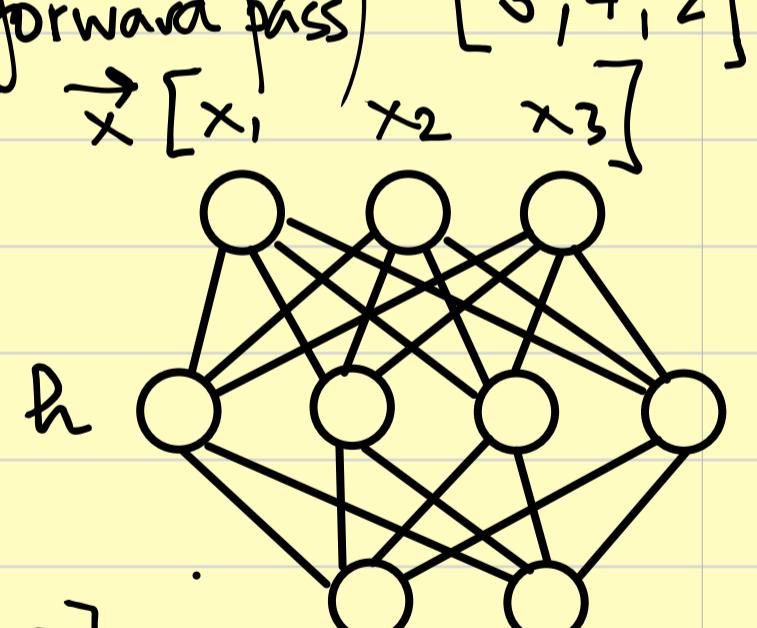
$$w_{h2} \quad b_2$$

$$\begin{bmatrix} -1 \\ 3 \\ 5 \\ 3 \end{bmatrix} \xrightarrow{\text{ReLU}} \begin{bmatrix} 0 \\ 3 \\ 5 \\ 3 \end{bmatrix} [1 \times 2]$$

$$\hat{y} [ \hat{y}_1 \quad \hat{y}_2 ]$$

$$\begin{bmatrix} 1 & 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 & 1 \end{bmatrix} \xrightarrow{\text{ReLU}} \begin{bmatrix} -2 \\ 3 \end{bmatrix}$$

$$\hat{y} [ \hat{y}_1 \quad \hat{y}_2 ]$$



$$f(\sum w_{h_2} \cdot \hat{h} + b_2) = f\left(\begin{bmatrix} 1 & 1 & -1 & 0 \\ 0 & 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 3 \\ 5 \\ 3 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix}\right) = f\left(\begin{bmatrix} -2 \\ 3 \end{bmatrix}\right) = \begin{bmatrix} 0 \\ 3 \end{bmatrix}$$

$1.0 + 1.3 + (-1) \cdot 5 + 0.3 + 0 = -2$   
 $0.0 + 0.3 + 1.5 + (-1) \cdot 3 + 1 = 3$

---

Example 4. Perception with hidden layer  $[3, 4, 2, 3]$  with feed forward & Back propagation.

SOFT MAX  $\vec{x} = \begin{bmatrix} 2 \\ 1 \\ 3 \end{bmatrix} \rightarrow \text{SoftMax}(\vec{x}) = \begin{bmatrix} \frac{2}{2+1+3} \\ \frac{1}{2+1+3} \\ \frac{3}{2+1+3} \end{bmatrix} = \begin{bmatrix} 2/6 \\ 1/6 \\ 3/6 \end{bmatrix} = \begin{bmatrix} 0.33 \\ 0.16 \\ 0.5 \end{bmatrix}$

Decision Layer:

$\vec{x}$

2	2	-4	4	-2
1	1	-2	2	-1
3	3	-6	6	-3
1	1	-2	2	-1

dot product  
 $w_{1,h_1}^*$   
 $b_1^*$

1	-1	1	-5	-1
1	1	0	0	3
0	1	1	1	5
1	0	1	-2	3

$\approx$  ReLU

1	-2	2	-1
1	-1		

$w_{1,h_2}^*$   $b_2$  error<sub>1</sub>

1	-1	1	0	0
0	1	-1	1	3

2	4
---	---

$\approx$  ReLU

2	4
---	---

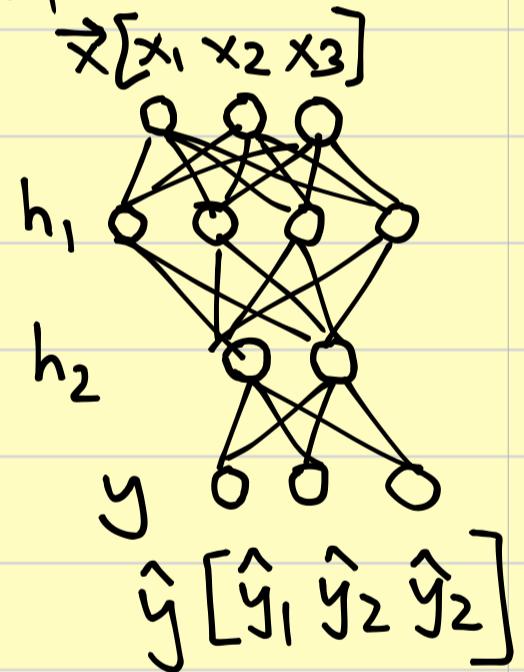
$w_{h_1, h_2}^*$   
dot product

$b_2^*$

$2 \cdot 0.23$   
 $2 \cdot (-0.77)$   
 $2 \cdot 0.54$  dot product!

②  $w_{h_2, 3}^*$

①  $b_3^*$



$0.23 \cdot 2 + 0 \cdot (-0.77) + 1 \cdot 0.54$   
 $0.23 \cdot 0 + 2 \cdot (-0.77) + 1 \cdot 0.54$

$w_{h_2, 3}^* b_3^*$

$\approx$  Soft Max (\*)

$\hat{y}$

$y$

$0.23 - 0 = 0.23$

$0.23 - (-0.77) = 1.00$

$0.23 - 0.54 = -0.31$

ERROR<sub>3</sub>

$$(*) \text{ SoftMax} \begin{bmatrix} 3 \\ 3 \\ 7 \end{bmatrix} = \begin{bmatrix} \frac{3}{3+3+7} \\ \frac{3}{3+3+7} \\ \frac{7}{3+3+7} \end{bmatrix} = \begin{bmatrix} 0.23 \\ 0.23 \\ 0.54 \end{bmatrix}$$

↑  
↑  
PREDICTION TARGET

Back Propagation:

- ① Transpose the error vector =  $\hat{y} - y$
- ② Dot Product of the transposed error  $\hat{y} - y$  and the output of previous layer
- ③ Multiply the transposed error  $\hat{y} - y$  with weight matrix of current layer.  $w_{h2,3} \cdot \text{error}_3^T$
- ④ Repeat until we are above.

Now we repeat our feed forward and backward propagation until the weights deliver a minimum in the cost function.

$$\begin{matrix} 2 \\ 1 \\ 3 \end{matrix}$$

$w_{1,h_1}^*$        $b_1^*$

2	1	3	1
-4	-2	-6	-2
4	2	6	2
2	-1	3	-1

$$\begin{matrix} 15 \\ -30 \\ 30 \\ 11 \end{matrix} \xrightarrow{\sim \text{ReLU}} \begin{matrix} 15 \\ 0 \\ 30 \\ 11 \end{matrix}$$

$$w_{h_1, h_2}^* \quad b_2^*$$

0	3	5	3	1
0	-3	-5	-3	-1

$$\begin{matrix} 184 \\ -184 \end{matrix} \xrightarrow{\sim \text{ReLU}} \begin{matrix} 184 \\ 0 \end{matrix}$$

0'46	0'92	0'23	84'87		0'15	0
-1'54	-3'08	-0'77	-284'13	$\approx$	0'5	1
1'08	2'16	0'54	199'26	Soft Max	0'35	0

$$\text{SoftMax} \begin{bmatrix} 84'87 \\ -284'13 \\ 199'26 \end{bmatrix} = \begin{bmatrix} 84'87 \\ 284'13 \\ 199'26 \end{bmatrix} \cdot \frac{1}{568'26} = \begin{bmatrix} 0'15 \\ 0'5 \\ 0'35 \end{bmatrix}$$

$= \begin{bmatrix} -0'15 \\ -0'5 \\ -0'35 \end{bmatrix} \text{error}_3^*$

$84'87 + 284'13 + 199'26 = 568'26$

