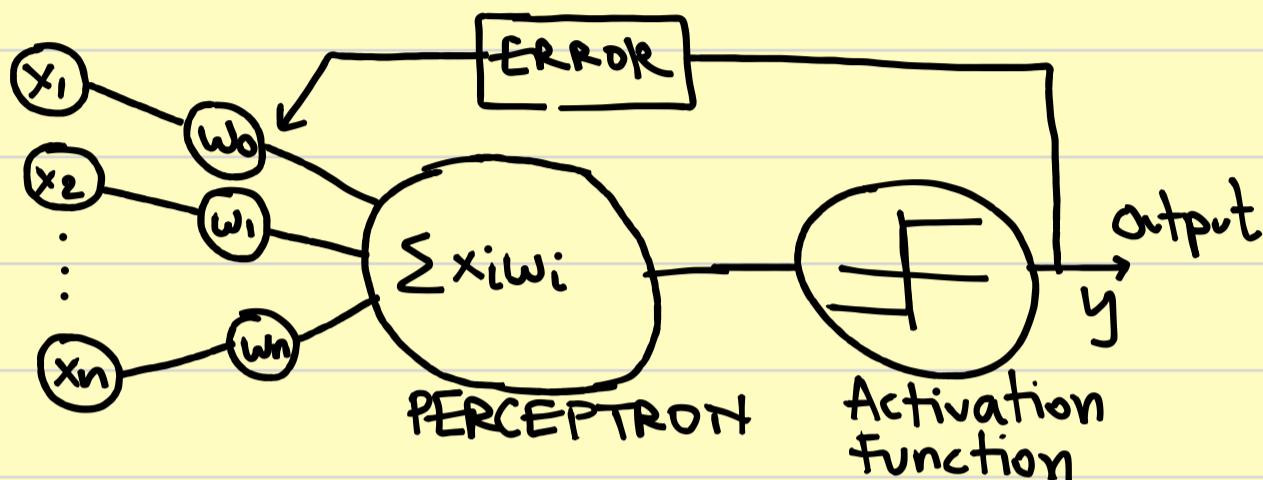


DEEP LEARNING by hand

SINGLE NEURON PERCEPTRON

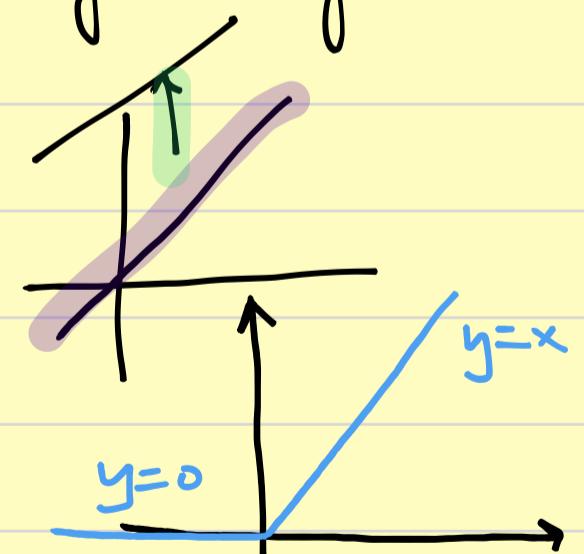
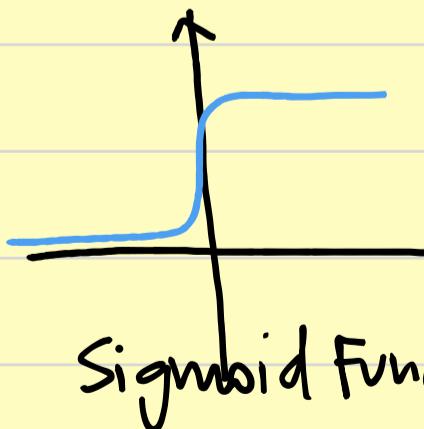
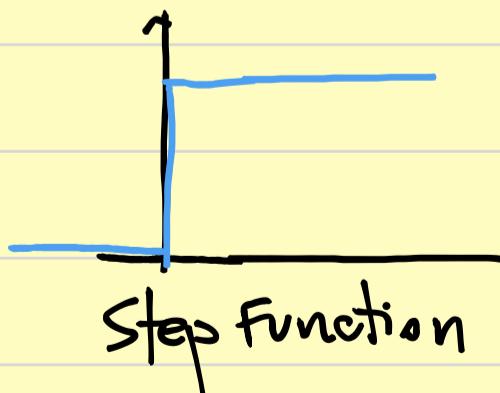


Step 1. Multiply all input values x_i with their corresponding weights w_i and then calculate the weighted sum.

$$\sum_{i=1}^n x_i w_i = x_1 w_1 + x_2 w_2 + \dots + x_n w_n$$

Step 2. An activation function is applied to the above giving us an output of the form

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



Rectifier
linear
Unit

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 y .

$$\text{We calculate the error: } C = \frac{1}{2} (y - \hat{y})^2 = \\ = \frac{1}{2} (y - \sum w_i x_i - b)^2$$

We call this error "cost" and we want to find the values of w_1, w_2, \dots, w_n that minimize the cost.

$$\frac{dC}{dw_1} = \frac{1}{2} \cdot 2 \cdot (y - w_1 x_1 - w_2 x_2 - \dots - w_n x_n) (-x_1)$$

$$\frac{dC}{dw_2} = \frac{1}{2} \cdot 2 \cdot (y - w_1 x_1 - w_2 x_2 - \dots - w_n x_n) (-x_2)$$

$$\vdots \\ \frac{dC}{dw_n} = \frac{1}{2} \cdot 2 \cdot (y - w_1 x_1 - w_2 x_2 - \dots - w_n x_n) (-x_n)$$

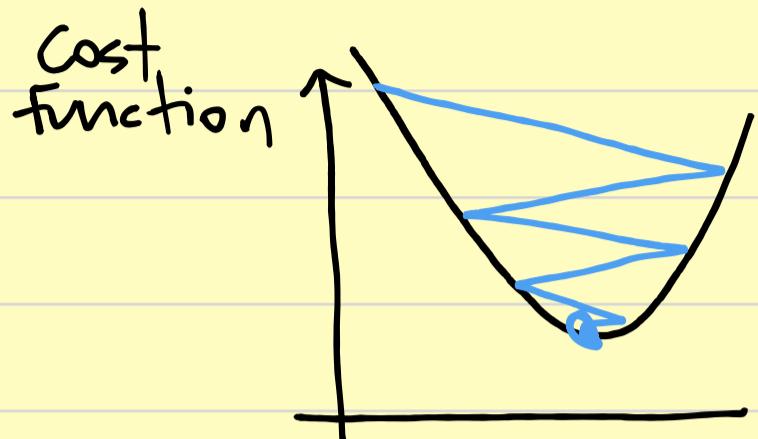
These derivatives are called gradients.

Step 4. Now to reach a minima we will start taking small steps towards the direction of the minima (opposite direction of the gradient).

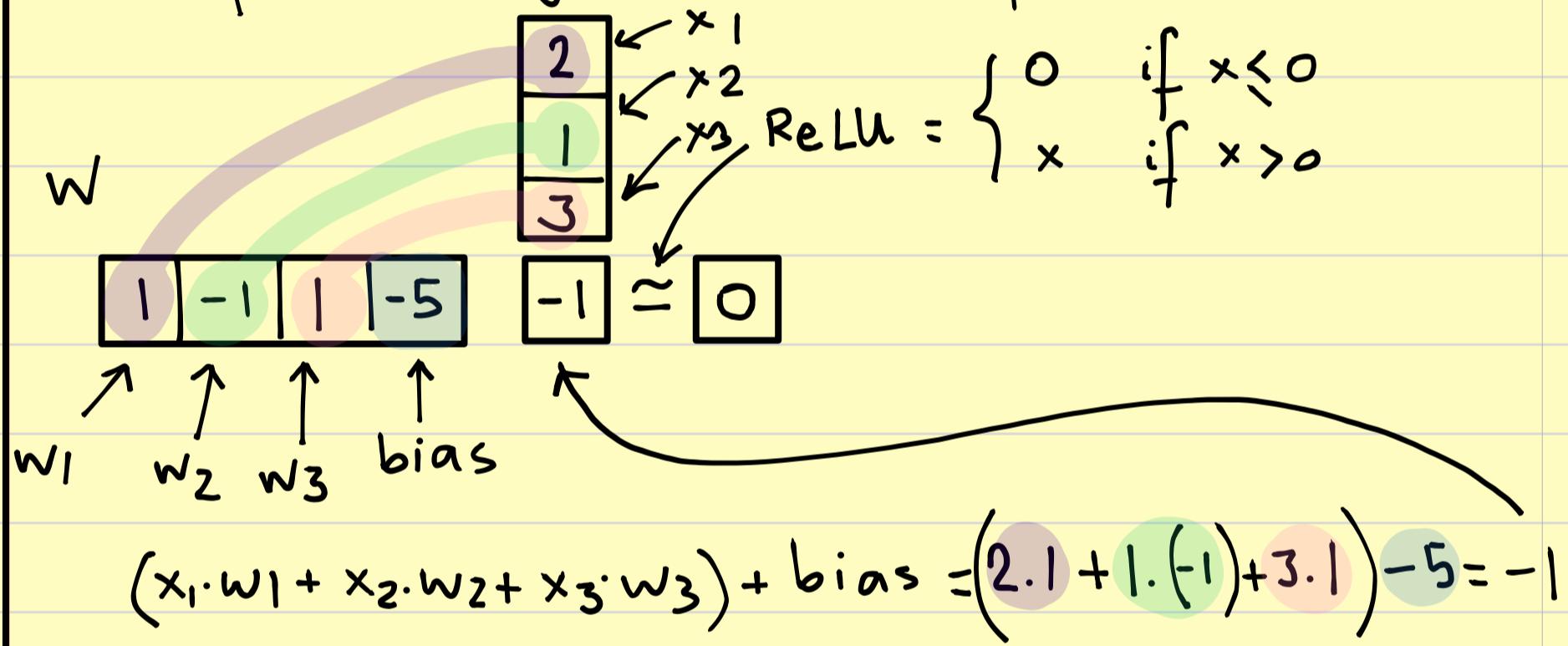
GRADIENT DESCENT.

We do so by updating the weights by a small

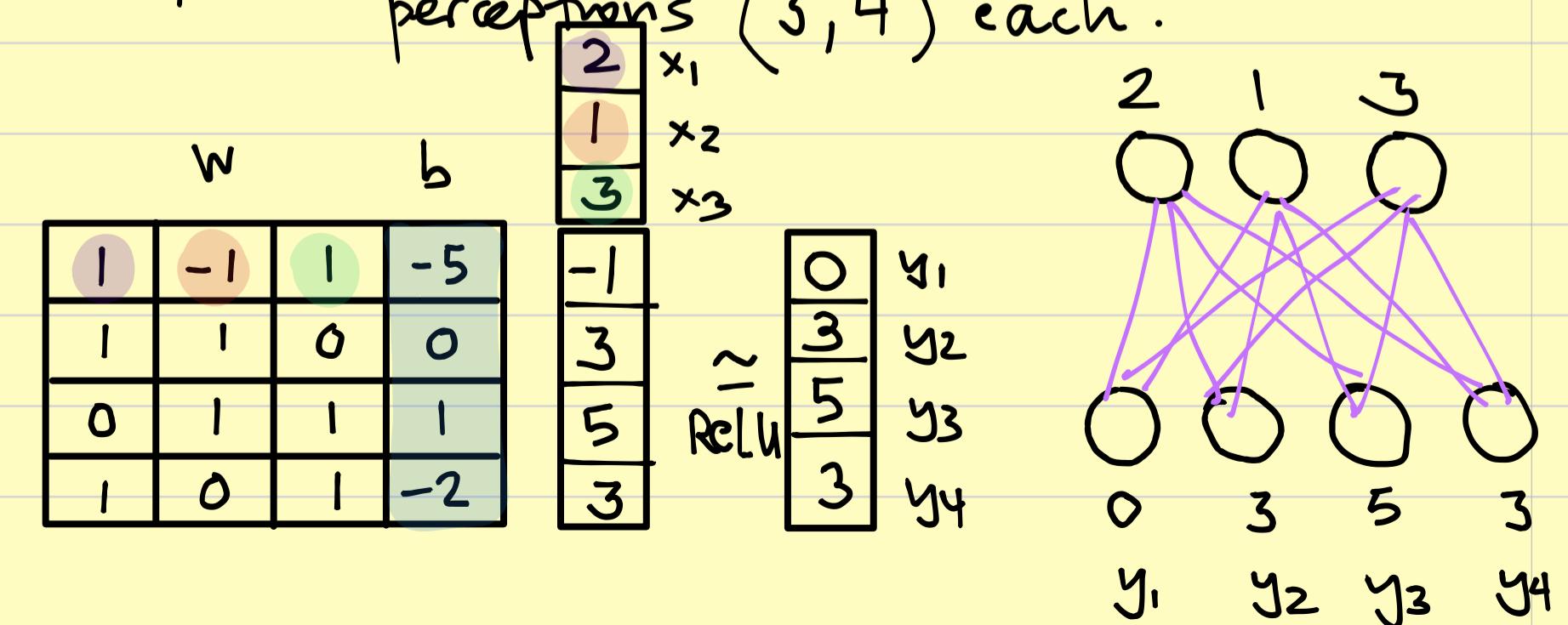
amount each time. This small amount is called **LEARNING RATE**.



Example 1. Single Neuron Perception Forward Pass

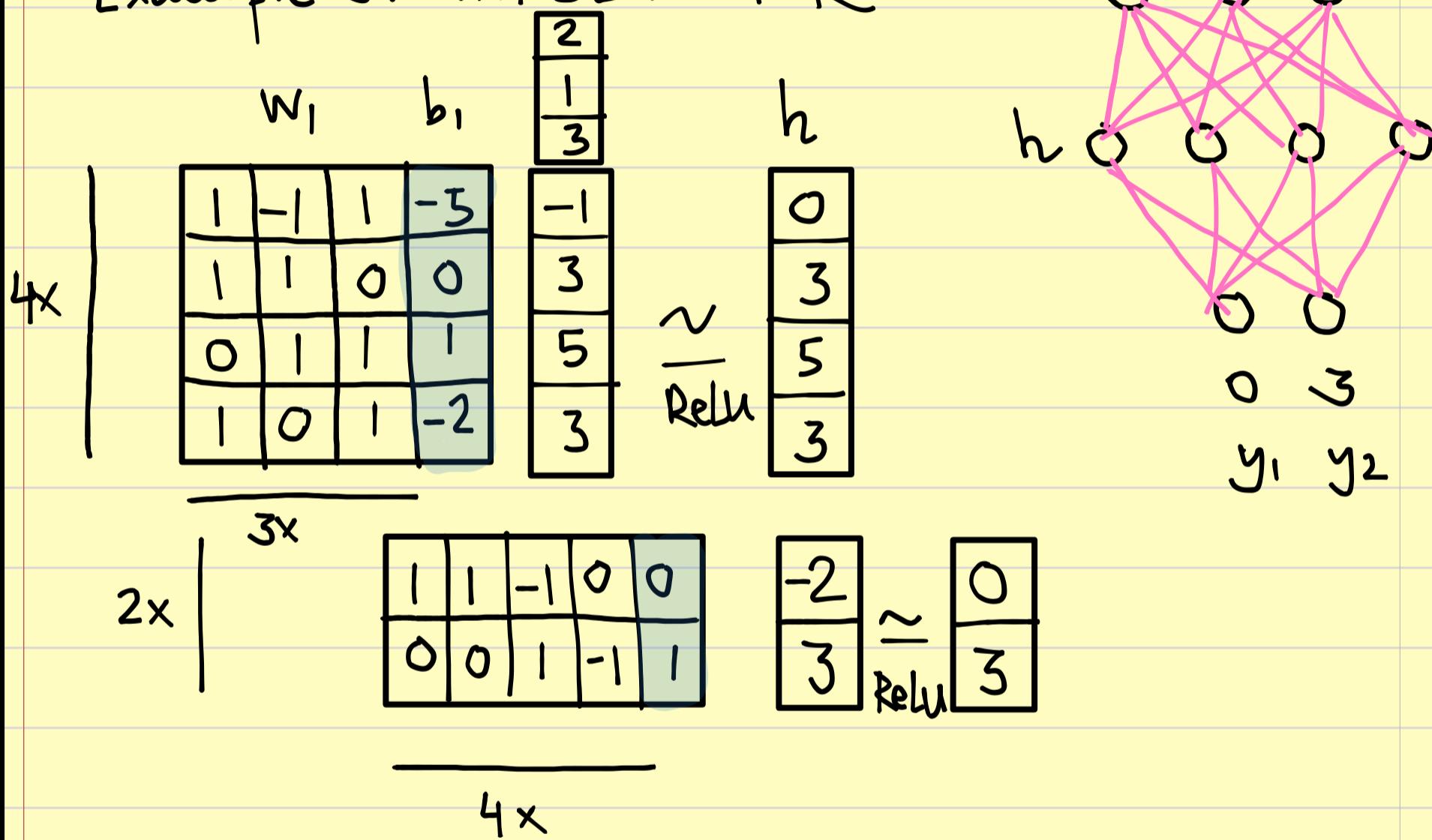


Example 2. Calculate the output of two layers of perceptions (3, 4) each.



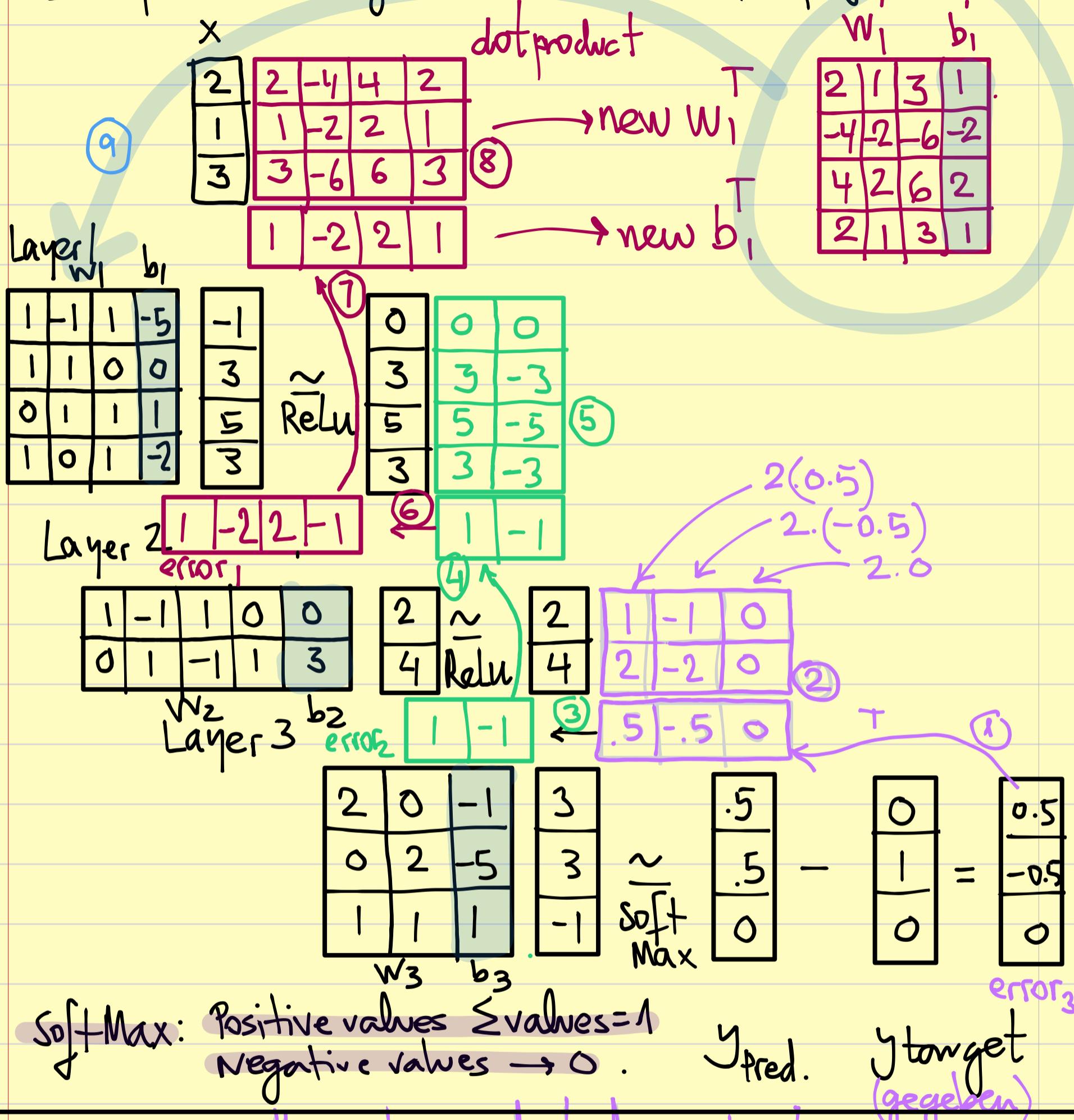
$$\begin{aligned}
 & (2.1 - 1.1 + 1.3) - 5 = -1 \\
 & (2.1 + 1.1 + 0.3) + 0 = 3 \\
 & (2.0 + 1.1 + 1.3) + 1 = 5 \\
 & (2.1 + 0.1 + 1.3) - 2 = 3
 \end{aligned}$$

Example 3. HIDDEN LAYER



$$\begin{aligned}
 & (1.0 + 1.3 - 1.5 + 0.3) + 0 = -2 \\
 & (0.0 + 0.3 + 1.5 - 1.3) + 1 = 3
 \end{aligned}$$

Example 4. Feed forward and Backpropagation.



① Transpose the vector predicted - vector target

② Multiply the vector predicted error₃ with

Dot prod of Layer 2. DOT. PRODUCT!

③ Multiply the predicted error₃ with w_3 matrix.

$$(0.5)2 + (-0.5)0 + 0.1 = 1$$

$$(0.5) \cdot 0 + (-0.5) \cdot 2 + 0.1 = -1$$

④ Transpose of the $\text{error}_3 \cdot W_3$

⑤ DOT PRODUCT of $(\text{error}_3 \cdot W_3)^T \cdot \text{output Layer 1}$.

⑥ Multiply the predicted error₂ with W₂ matrix.

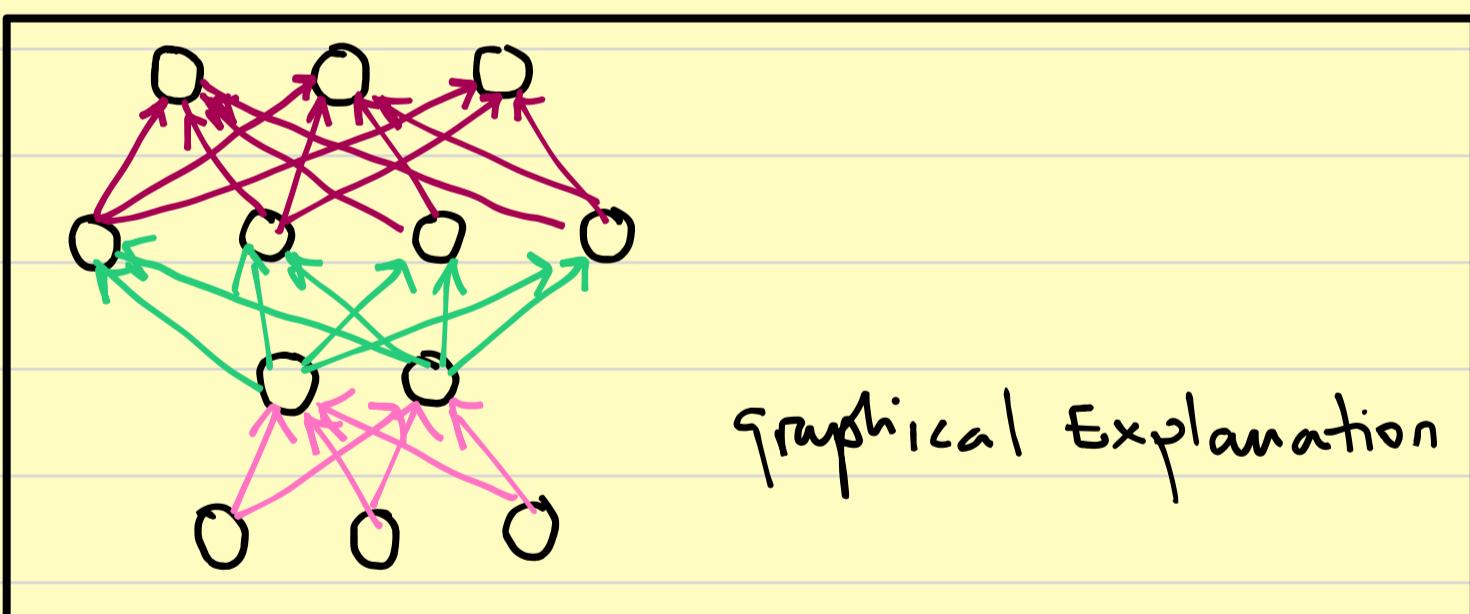
$$(1.1) + (-1) \cdot 0 = 1$$

$$(1)(-1) + (-1) \cdot 1 = -2$$

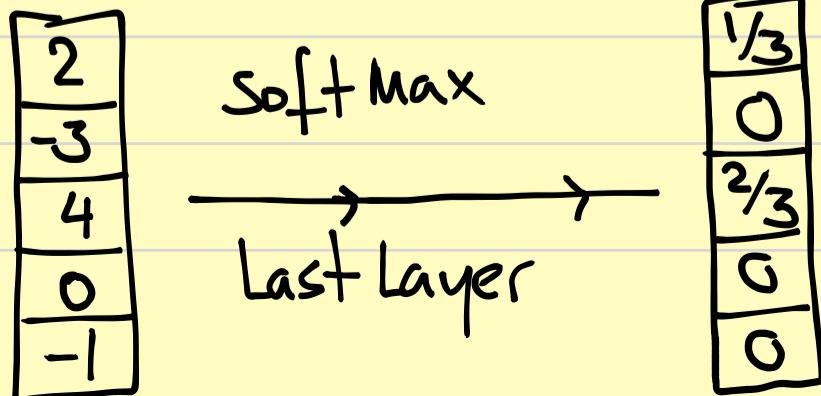
⑦ Transpose of the $\text{error}_2 \cdot W_2$

⑧ DOT PRODUCT of $(\text{error}_2 \cdot W_2)^T \cdot \text{Input}$

⑨ Substitute the new updated weights w'_1 & b'_1



Softmax explanation: ≤ 1



$$2 + 4 = 6$$

$$\frac{2}{6} + \frac{4}{6} = 1$$

$$\frac{1}{3} + \frac{2}{3}$$

TRANSFORMERS

Transformers are neural networks that eliminate the need for labeled datasets and convolution in the network.

The two main mechanisms that allow for this are:

1) ATTENTION WEIGHTING

2) FEED FORWARD NETWORKS (FFN)

Attention weighting is a technique by which the model learns which part of the incoming sequence needs to be focused on. It screens the data and recognizes what is relevant.

FFN are in this context a regular multilayer perception acting on a batch of data vectors.

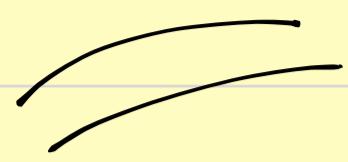
Combined FFN & AW, it produces "position dimension".

FFN mechanism (see above ML Perception).

What is the attention mechanism?

The transformer architecture can be seen as a

combination of two shells :



The outer shell is a combination of attention-weighting and a FFN.

The inner shell consists on a self-attention mechanism and is part of the attention-weighting feature.

Example : The attention weight matrix A (given) is obtained by feeding the input features into the architecture. This matrix tries to find the most relevant parts in the input sequence.

x_1	x_2	x_3	x_4			
2	0	0	2			
0	1	0	0			
0	2	1	0			
0	0	1	1			
2	0	0	0			
1	0	1	1			

① Weight W_Q

1	1	0	0	0	0
0	1	0	1	0	0
0	0	1	0	1	1

2	1	0	2
0	1	1	1
3	2	2	1

$$W_Q \cdot x = q$$

QUERY
VECTOR

② key

1	0	0	0	1	0	0
0	1	0	0	0	0	0
3	0	0	1	0	0	0

$$W_K \cdot x = k$$

One-hot-encoding

$$W_V \cdot x = v$$

③ Value W_V

④ The next step is to multiply $k \cdot q$ with dot product

5 The next step is to scale/normalize each element

$$z_i = \frac{x_i - \mu}{\sigma_x}$$

6 The next step is to apply softmax to get the prediction

7 Backwards propagation as in FFN.

Apfel → 000

Mango → 011

Orange → 001

Ananas → 101

Erdbeer → 010

Zitrone → 110

Birne → 100

Melone → 111

$$\begin{bmatrix} 0.01 \\ 0.99 \\ 0.0 \end{bmatrix} \rightarrow \text{Erdbeer}$$

y_{pred}

