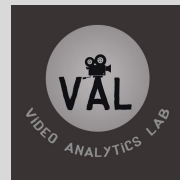


LEARNING WITH NEURAL NETWORKS



By Nithish Divakar
Video Analytics Lab
Dept. Computational and Data Sciences
IISc Bangalore



What is the following number ?

504192

INSPIRATION

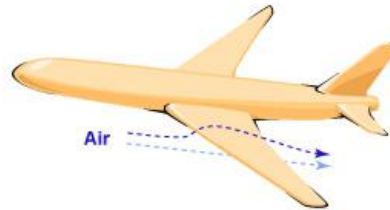
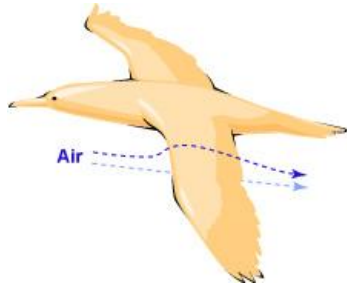
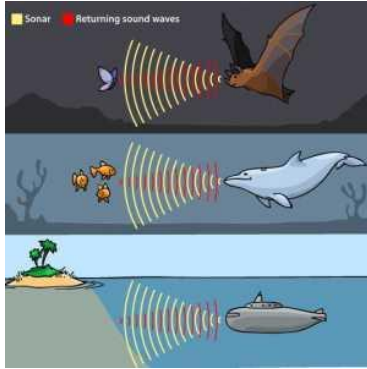
What is the following number ?

504192

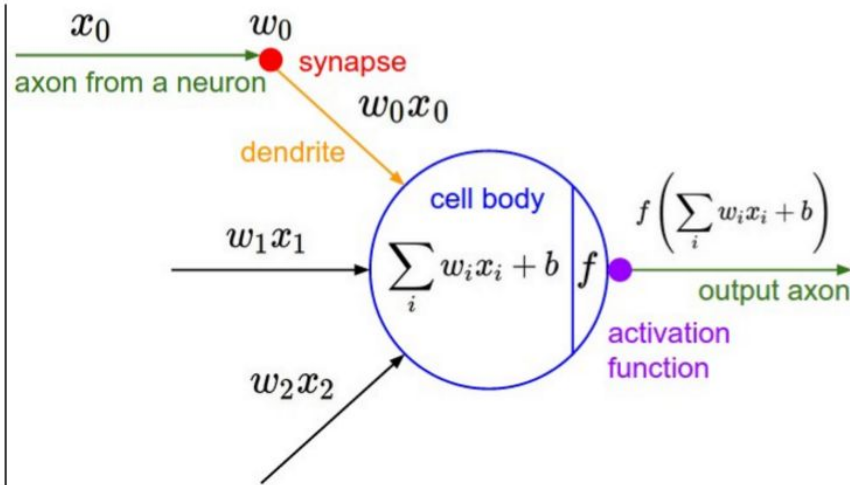
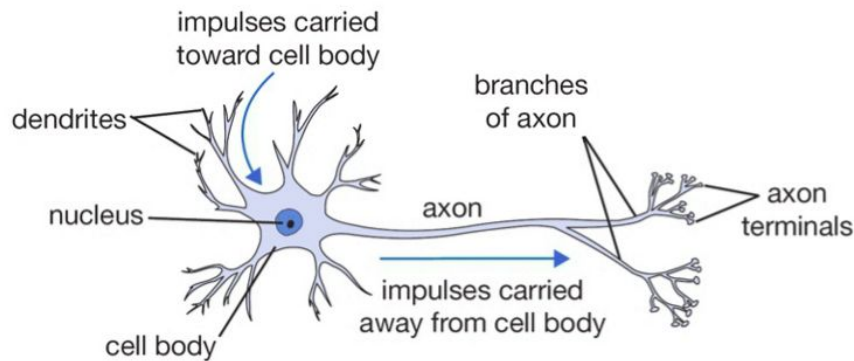
WHY NOT COPY THE BEST 'NATURAL' LEARNING SYSTEM ...

... TO CREATE THE BEST 'ARTIFICIAL' LEARNING SYSTEM ?

COPYING FROM NATURE -- NOTHING NEW, REALLY !!!



BASIC BUILDING BLOCK :: A NEURON

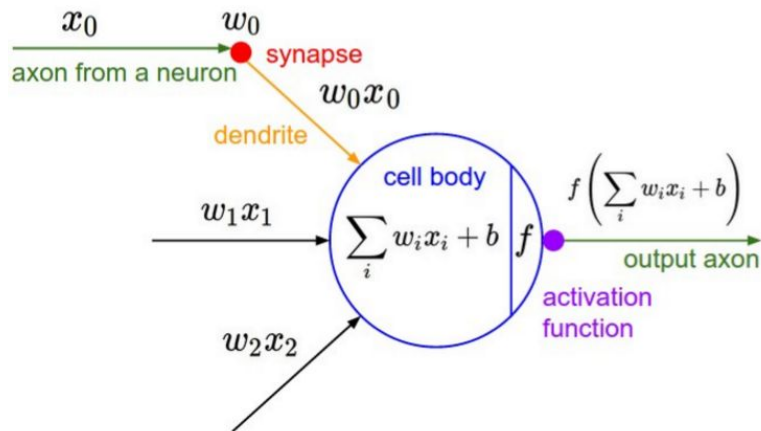


MAPPING TO OUR PROBLEM

$$y = f(\mathbf{x}; \mathbf{w})$$

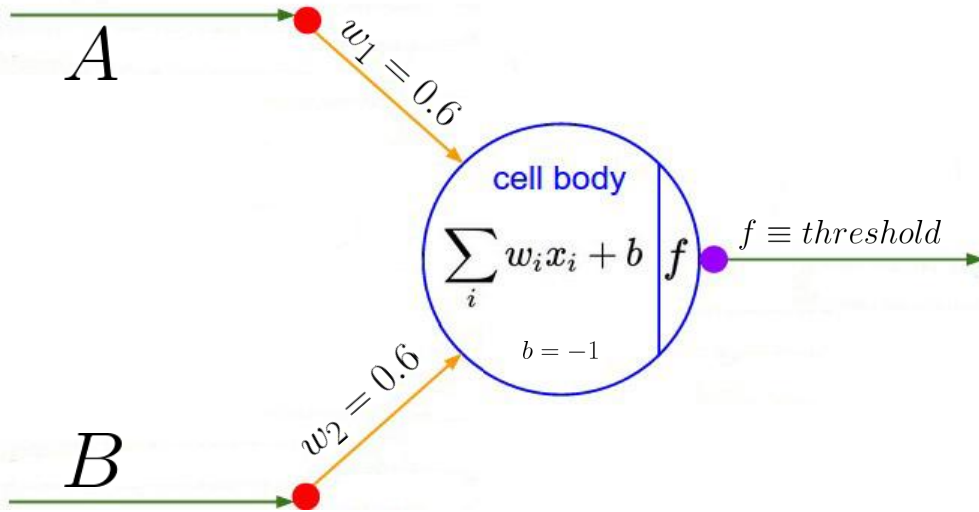
– \mathbf{x}, y : known

– \mathbf{w} : weights (need to be “learnt”)

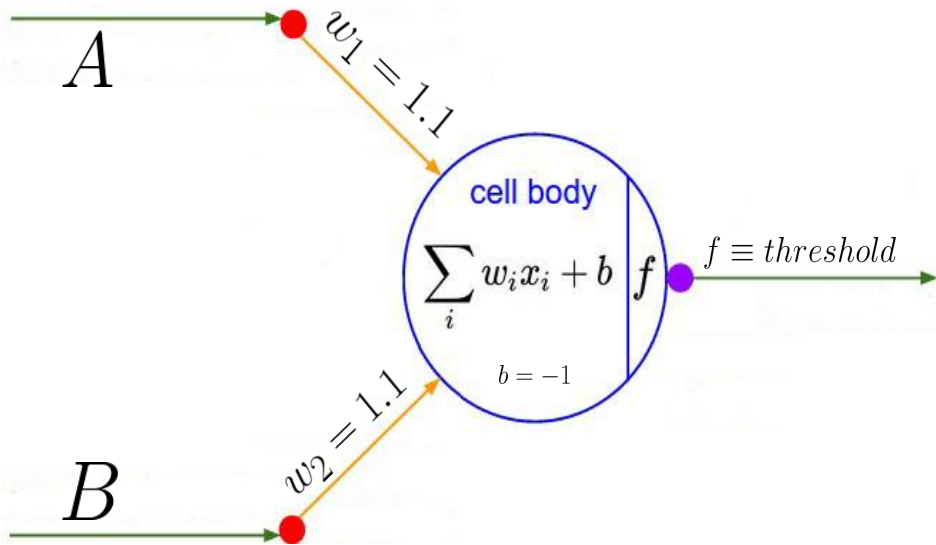




LOGICAL FUNCTIONS ... FIRST NN IMPLEMENTATIONS (1943)



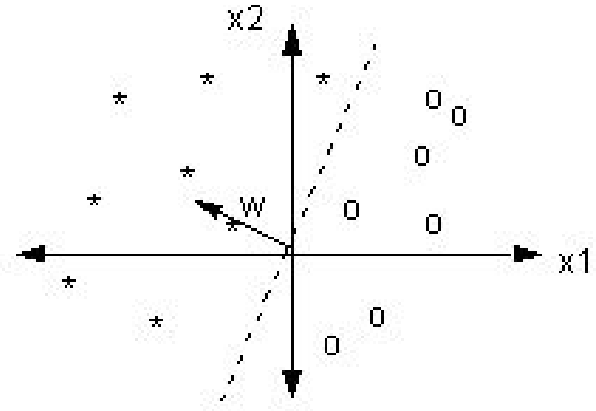
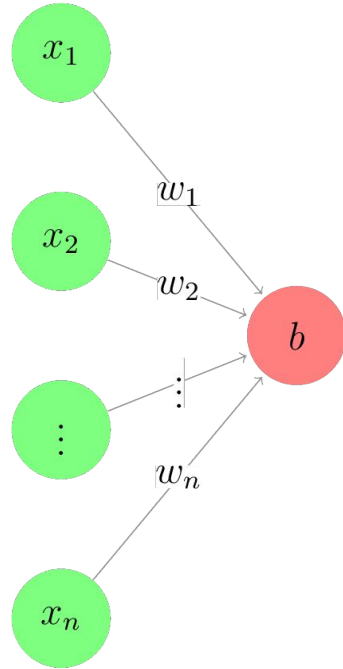
A	B	A and B
0	0	0
0	1	0
1	0	0
1	1	1



A	B	A or B
0	0	0
0	1	1
1	0	1
1	1	1

NOT GATE ?

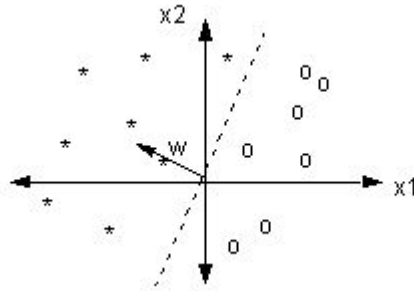
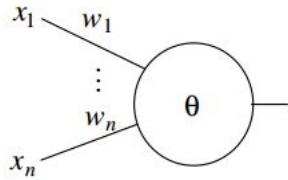
ROSENBLATT'S "PERCEPTRON" (1953)



LINEARLY SEPARABLE
CLASSES

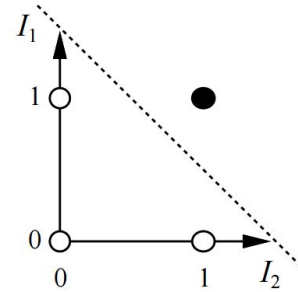
$$\text{output} = \begin{cases} 0 & \sum_i w_i x_i + b \leq 0 \\ 1 & \sum_i w_i x_i + b > 0 \end{cases}$$

ROSENBLATT'S "PERCEPTRON" (1953)



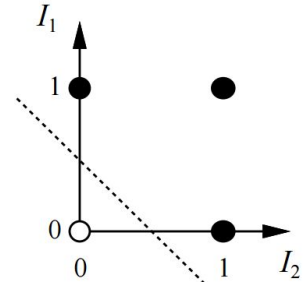
LINEARLY SEPARABLE
CLASSES

AND



(a) I_1 and I_2

OR



(b) I_1 or I_2

PERCEPTRON ALGORITHM

Initialize weights \mathbf{w} ; // Random or 0

```
Until [all examples correctly classified]
  For each training sample ( $\mathbf{x}, y$ )
    Compute  $y_t := \mathbf{x}^T \mathbf{w}$ 
    if  $y == y_t$  // Correctly classified
      continue ;
    else // Update weights
       $d\mathbf{w} = (y - y_t) * \mathbf{x}$  ;
       $\mathbf{w} = \mathbf{w} + d\mathbf{w}$ ;
    EndIf
  EndFor
EndUntil
```

PERCEPTRON ALGORITHM

Initialize weights w ; // Random or 0

Until [all examples correctly classified]

For each training sample (x, y)

Compute $y_t := x^T w$

if $y == y_t$ // Correctly classified
continue ;

else // Update weights

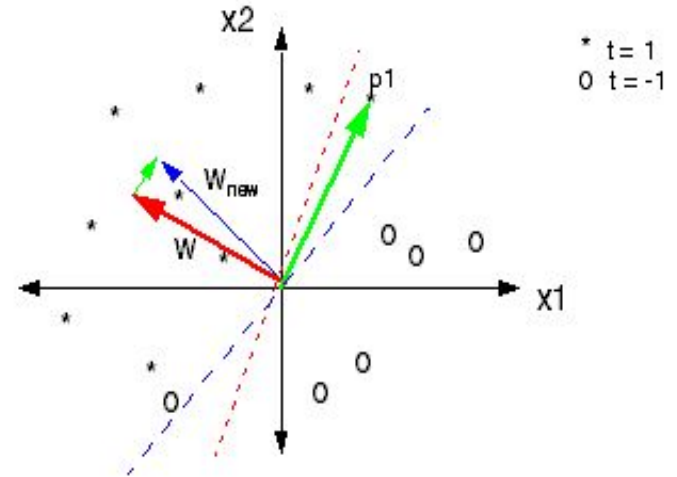
$dw = (y - y_t) * x$;

$w = w + dw$;

EndIf

EndFor

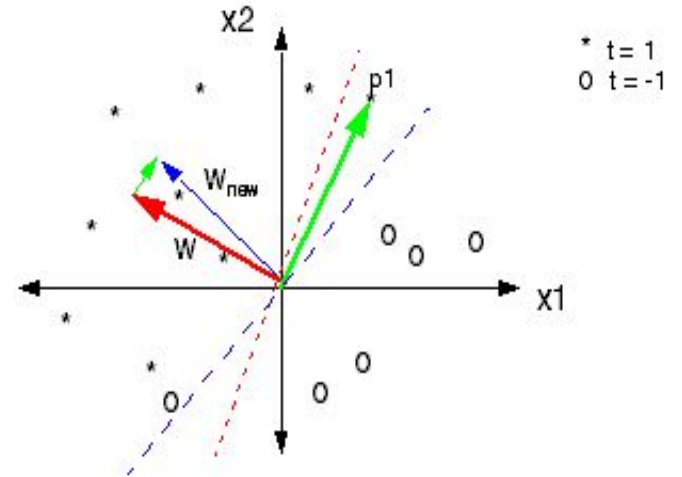
EndUntil



PERCEPTRON ALGORITHM

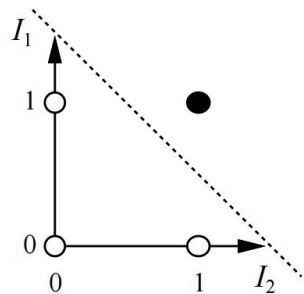
Demo ...

[visualize_PLA.html](#)



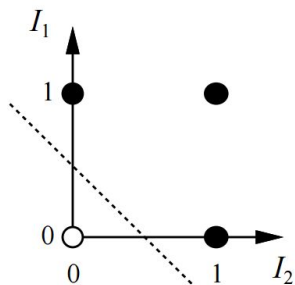
CAN A PERCEPTRON SEPARATE DISTRIBUTIONS WHICH ARE NOT LINEARLY SEPARABLE ???

AND

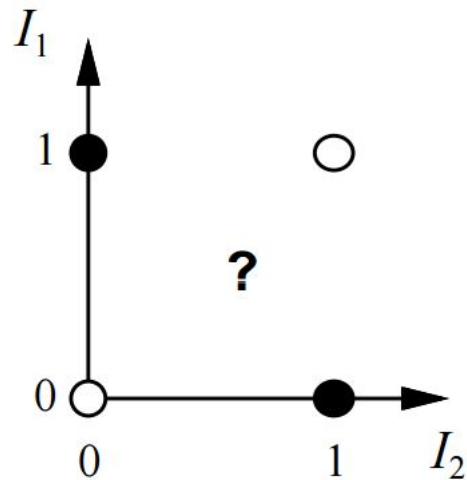


(a) I_1 and I_2

OR



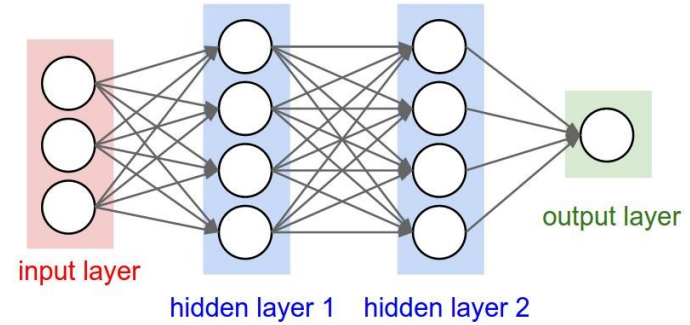
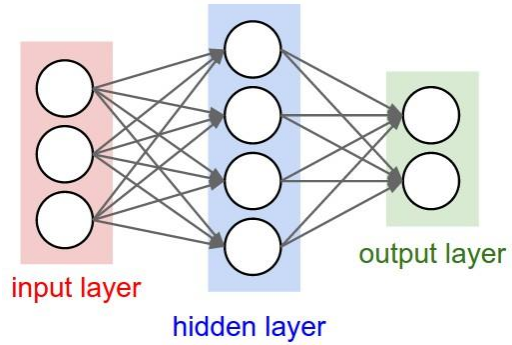
(b) I_1 or I_2



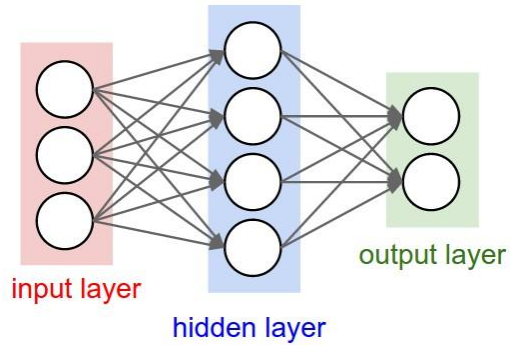
(c) I_1 xor I_2

LINEARLY SEPARABLE
CLASSES

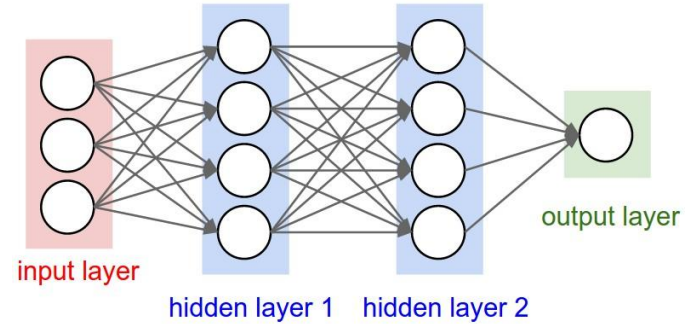
WHY USE ONLY ONE NEURON ?



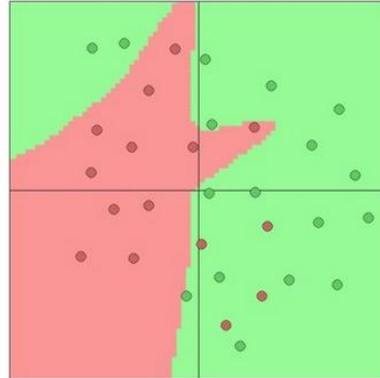
WHY USE ONLY ONE NEURON ?



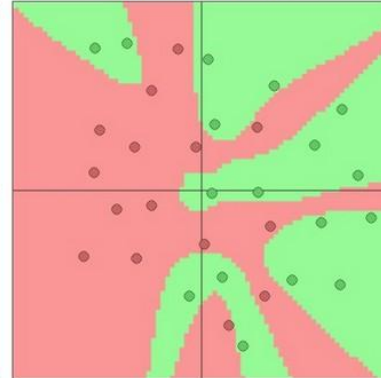
3 hidden neurons



6 hidden neurons



20 hidden neurons



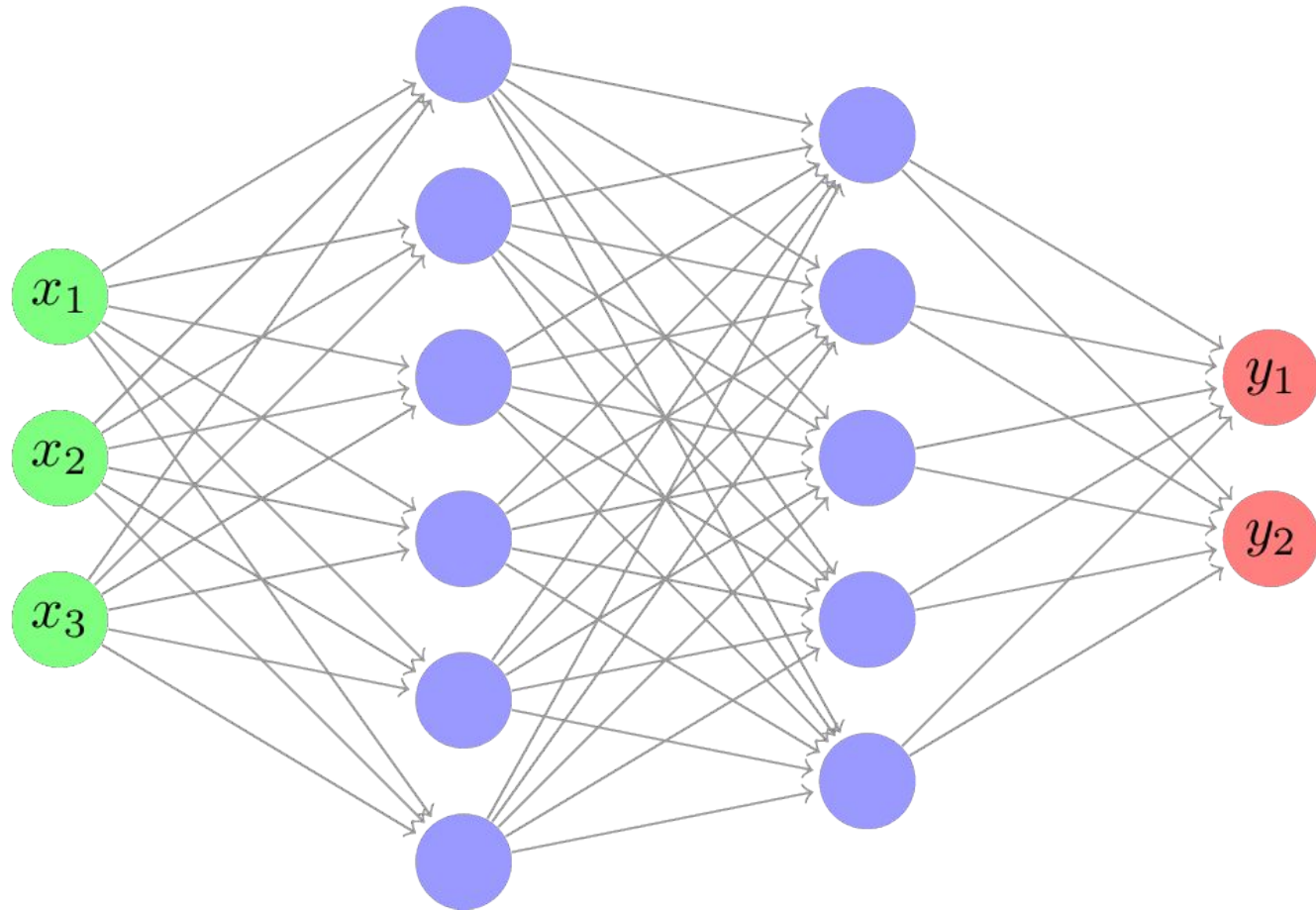
NEURAL NETWORKS

Input layer

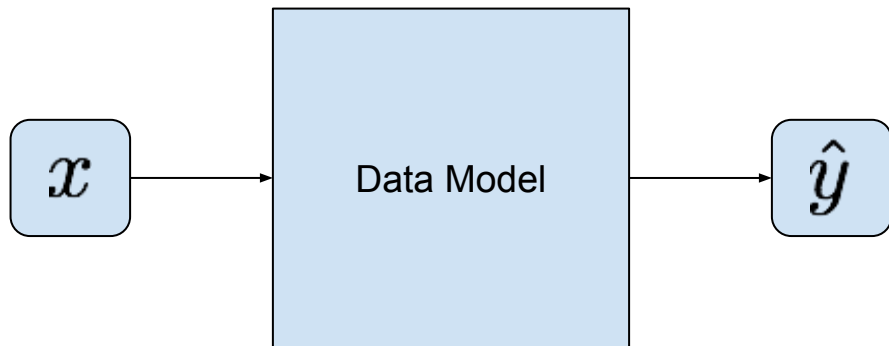
Hidden layer 1

Hidden layer 2

Output layer



DATA MODEL



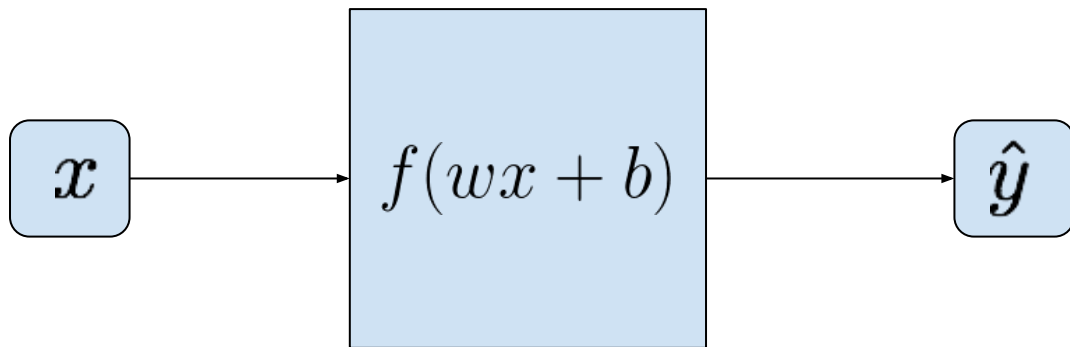
$$y \equiv \hat{y}$$

GOOD MODEL

$$y \neq \hat{y}$$

BAD MODEL

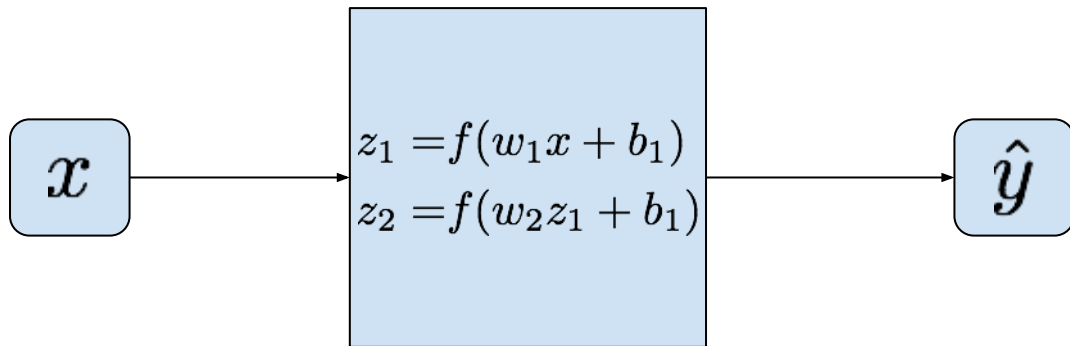
A SIMPLE DATA MODEL



$\{w, b\}$ parameters of the model

DEMO: [fitting_point.html](#)

A ~~SIMPLE~~ COMPLICATED DATA MODEL



$\{w_1, w_2, b_1, b_2\}$ parameters of the model

DEMO: [fitting_point2.html](#)

WHY DO WE NEED NO LINEARITY ?

$$z_1 = f(z_0)$$

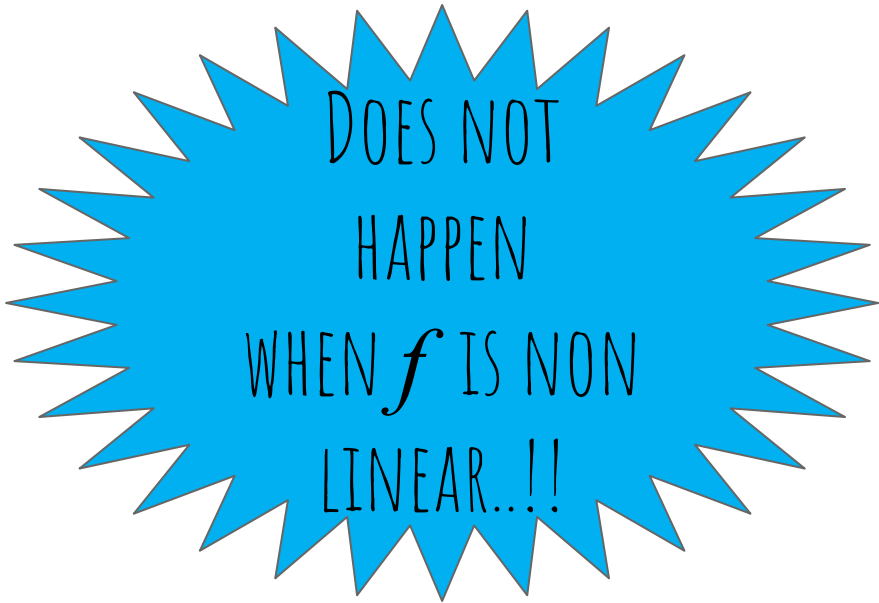
$$z_2 = f(z_1)$$

But if f is linear

$$y = f(x) = ax$$

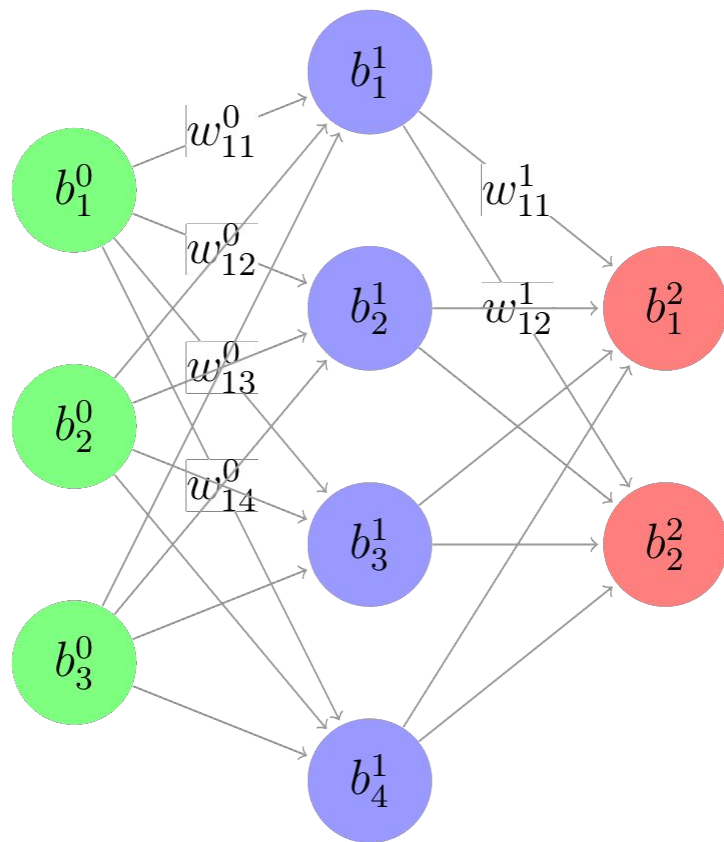
$$\begin{aligned}\implies z_2 &= f(z_1) \\ &= f(f(z_0))\end{aligned}$$

$$z_2 = f(az_0)..!!$$



DOES NOT
HAPPEN
WHEN f IS NON
LINEAR..!!

layer 0 layer 1 layer 2



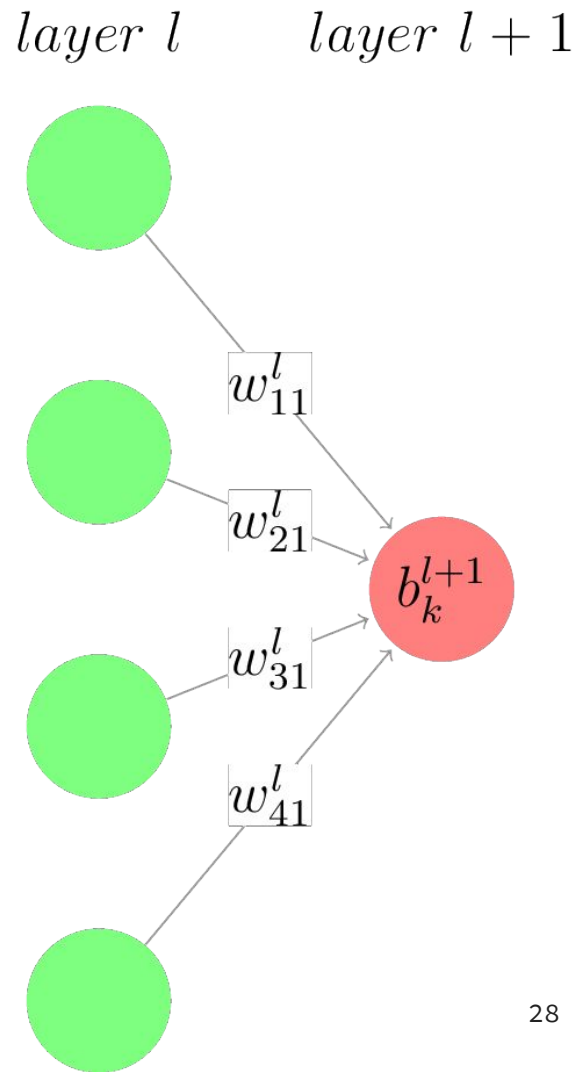
THINGS OF THE NETWORK

$w_{ji}^l \mapsto$ weight btw
 $node_j^l$ and $node_i^{l+1}$
 $b_{lj} \mapsto$ bias of $node_j^l$

OUTPUT OF A SINGLE NEURON

$$z_k^{l+1} = f \left(\sum_i w_{ji}^l z_i^l + b_i^{l+1} \right)$$

f IS THE NON LINEARITY
IN PERCEPTRON IT WAS THRESHOLD
FUNCTION



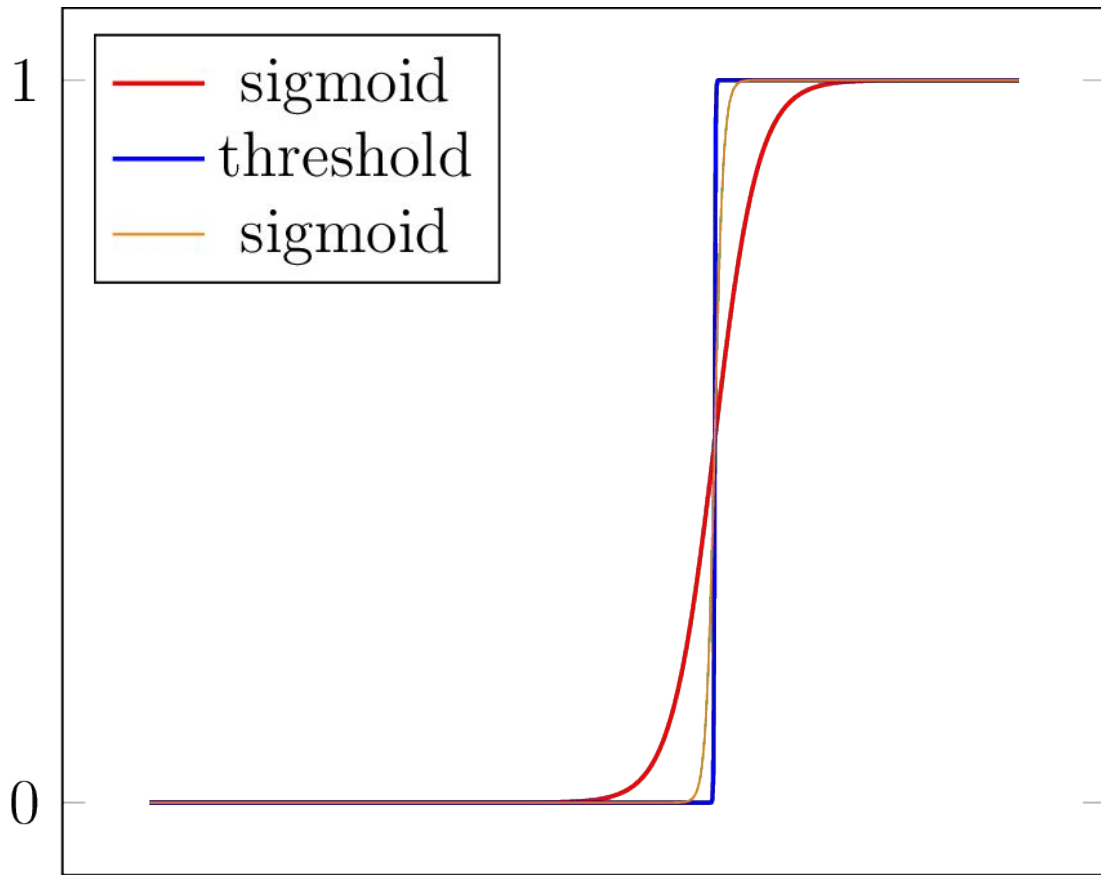
NON LINEARITY

$$z_k^{l+1} = f \left(\sum_i w_{ji}^l z_i^l + b_i^{l+1} \right)$$

$$w \leftarrow w + \nabla w$$

$$z = ?$$

$$\text{simoid}(z) = \frac{1}{1 + e^{-z}}$$



OUTPUT OF A LAYER

$$z_k^{l+1} = f \left(\sum_i w_{ji}^l z_i^l + b_i^{l+1} \right)$$

$$z_1^{l+1}$$

$$\sum_i w_{1i}^l z_i^l + b_1^{l+1}$$

$$z_2^{l+1}$$

$$\sum_i w_{2i}^l z_i^l + b_2^{l+1}$$

$$z_3^{l+1}$$

$$\sum_i w_{3i}^l z_i^l + b_3^{l+1}$$

⋮

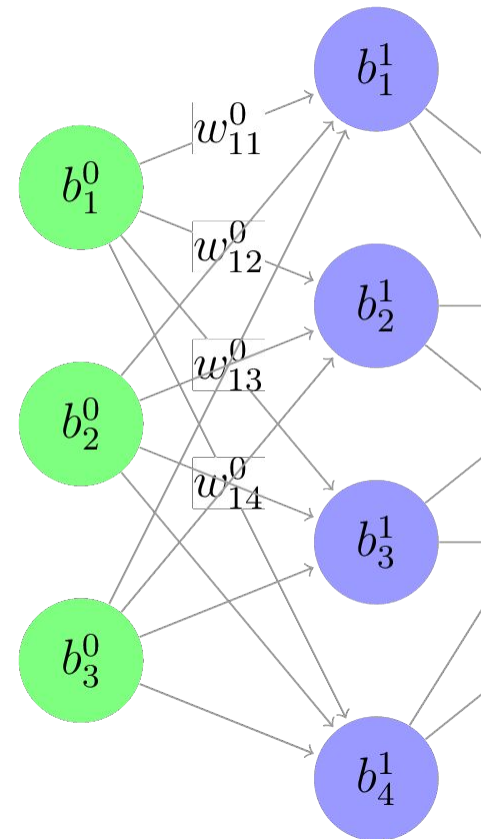
⋮

$$z_n^{l+1}$$

$$\sum_i w_{ni}^l z_i^l + b_n^{l+1}$$

layer 0

layer 1



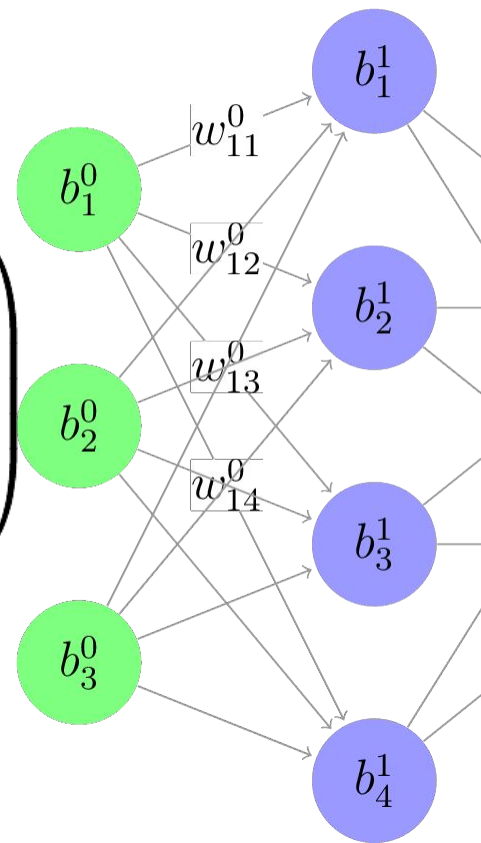
OUTPUT OF A LAYER

$$z_k^{l+1} = f \left(\sum_i w_{ji}^l z_i^l + b_i^{l+1} \right)$$

$$\begin{bmatrix} z_1^{l+1} \\ z_2^{l+1} \\ \vdots \\ z_m^{l+1} \end{bmatrix} = f \left(\begin{bmatrix} w_{11}z_1 + w_{12}z_2 \dots w_{1m}z_m \\ w_{21}z_1 + w_{22}z_2 \dots w_{2m}z_m \\ \vdots \\ w_{n1}z_1 + w_{n2}z_2 \dots w_{nm}z_m \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} \right)$$

layer 0

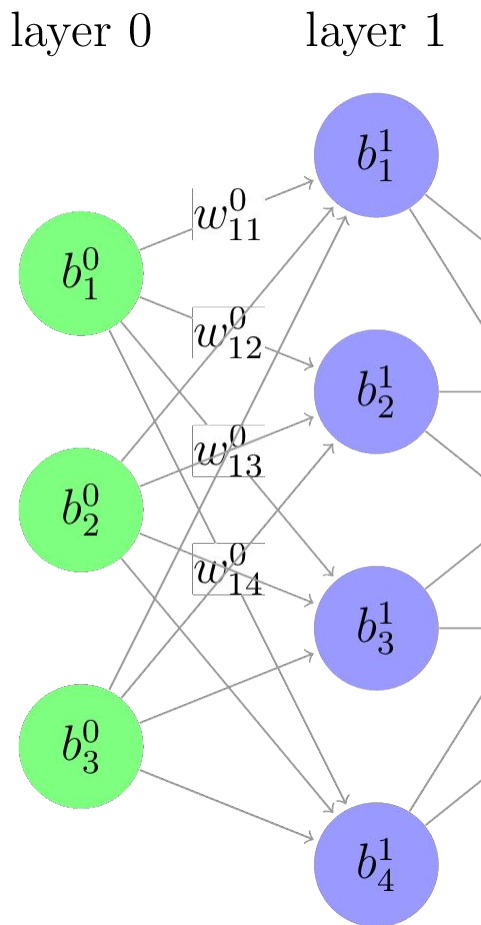
layer 1



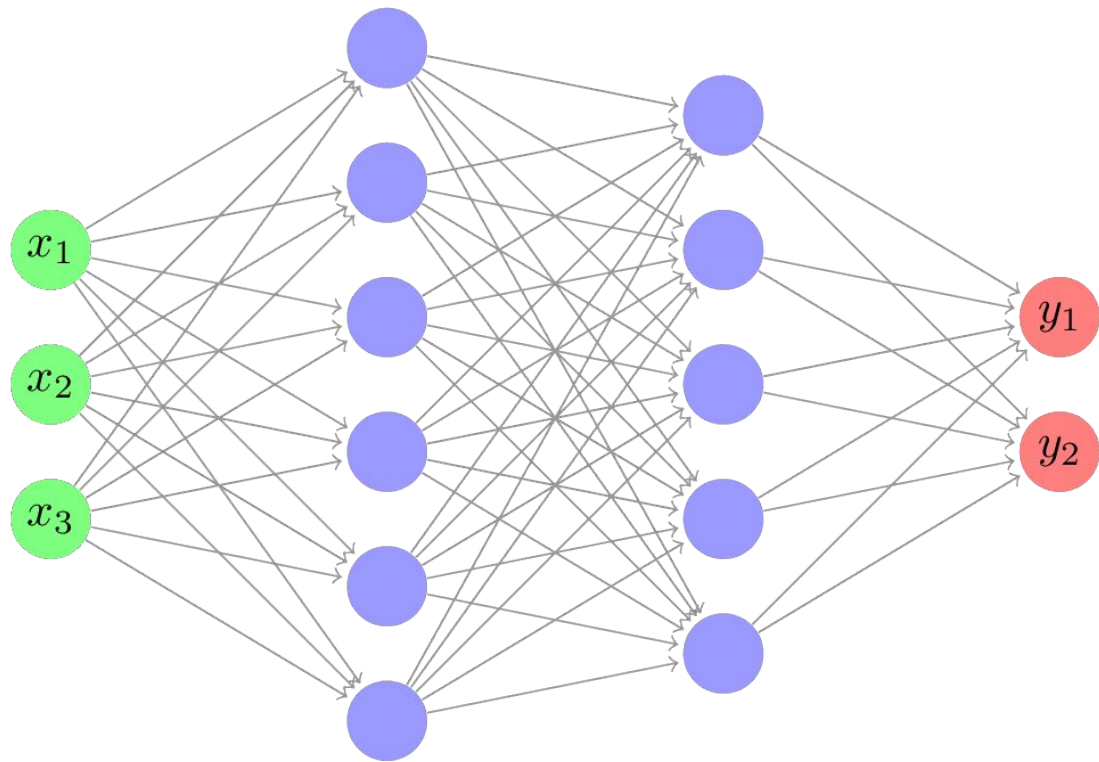
OUTPUT OF A LAYER

$$z_k^{l+1} = f \left(\sum_i w_{ji}^l z_i^l + b_i^{l+1} \right)$$

$$z^{l+1} = f(W^l z^l + b^{l+1})$$



NEURAL NETWORK



$$\hat{y} = f(f(f(f(f(W^1 x + b^2))))))$$

LEARNING THINGS!!

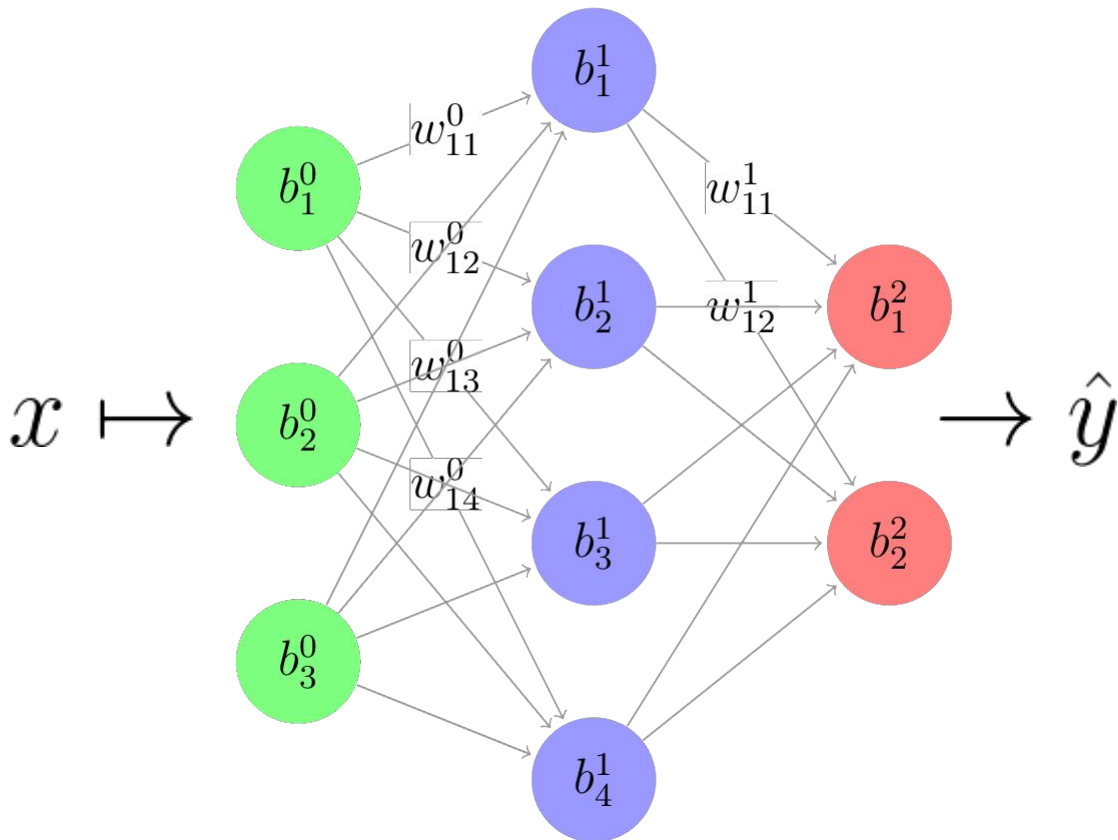
JUST ANOTHER NAME FOR FINDING
THE RIGHT PARAMETERS

?

layer 0

layer 1

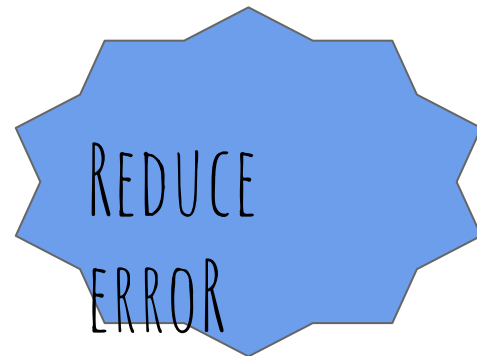
layer 2



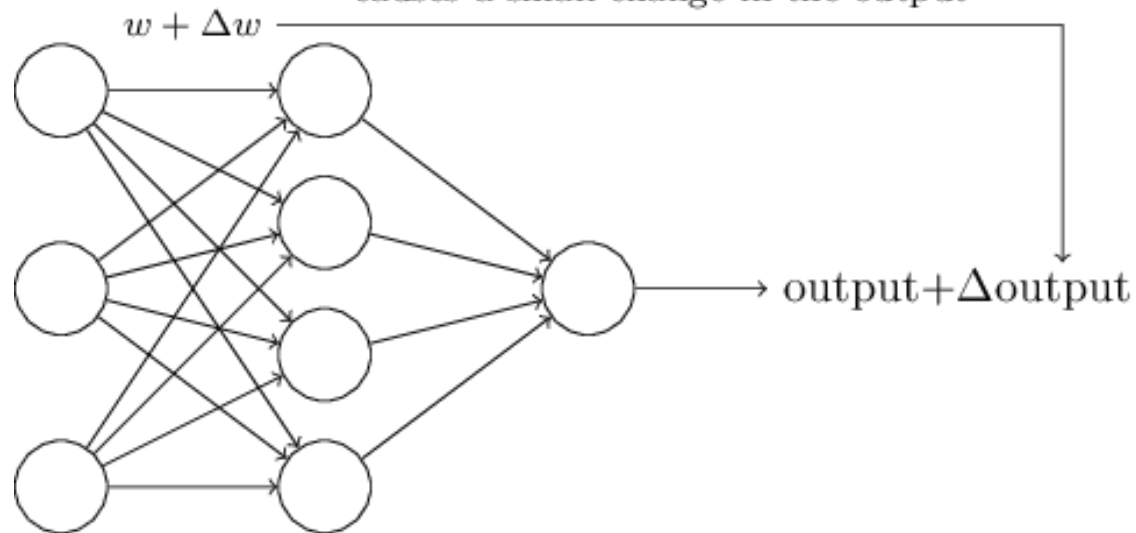
$$\text{er}[y, \hat{y}]$$

error

$$(y - \hat{y})^2$$



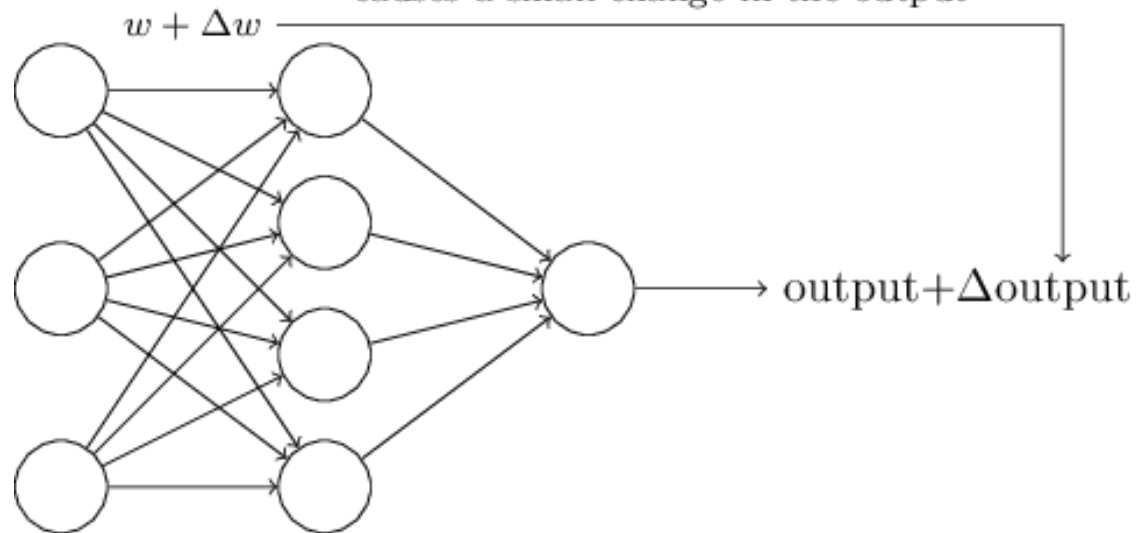
small change in any weight (or bias)
causes a small change in the output



WHAT HAPPENS
TO ERROR ?

$$(y - \hat{y})^2$$

small change in any weight (or bias)
causes a small change in the output

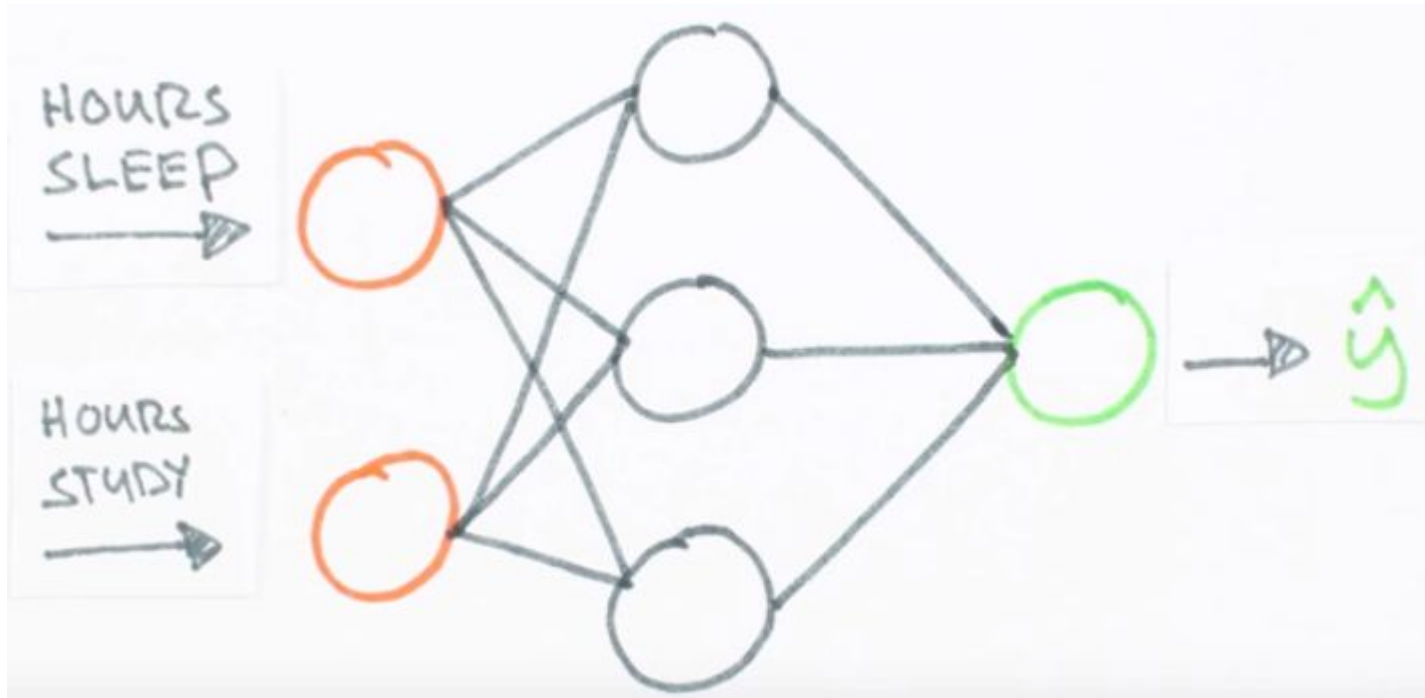


ANY WAY TO AUTOMATE THE PROCESS ? $(y - \hat{y})^2$

NEURON NETWORKS TRAINING

	X (HOURS SLEEP, HOURS STUDY)	y (SCORE ON TEST)
TRAINING	(3, 5)	75
	(5, 1)	82
	(10, 2)	93
TESTING	(8, 3)	?

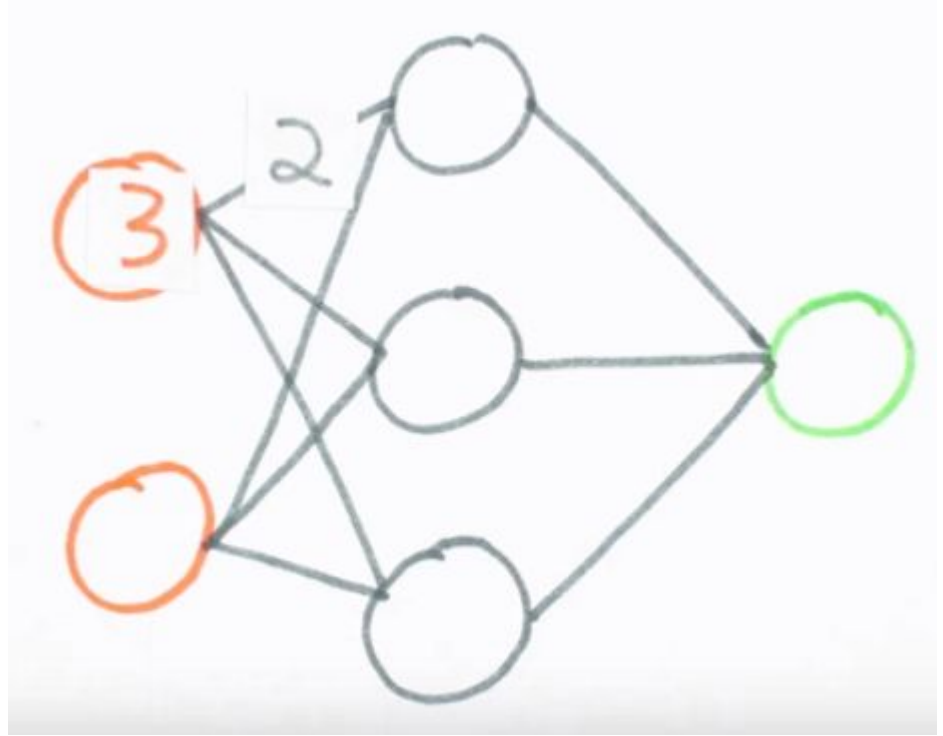
MULTI-NEURON NETWORKS :: ARCHITECTURE



MULTI-NEURON NETWORKS :: ARCHITECTURE

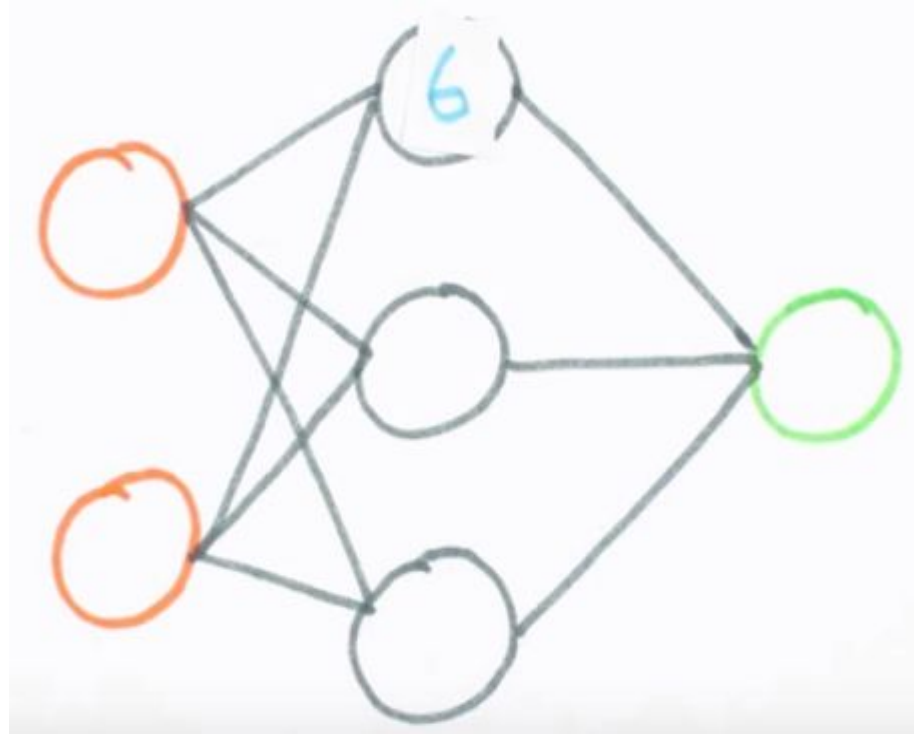
SYNAPSE

WEIGHT ON SYNAPSE =
2



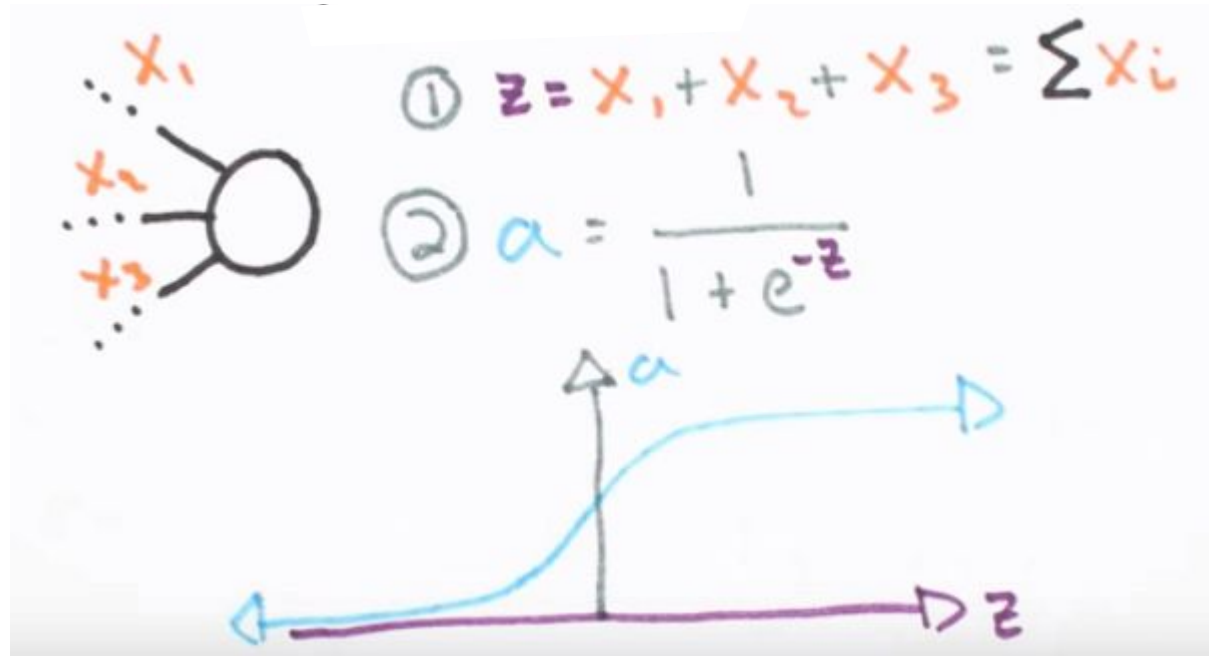
MULTI-NEURON NETWORKS :: ARCHITECTURE

SYNAPSE

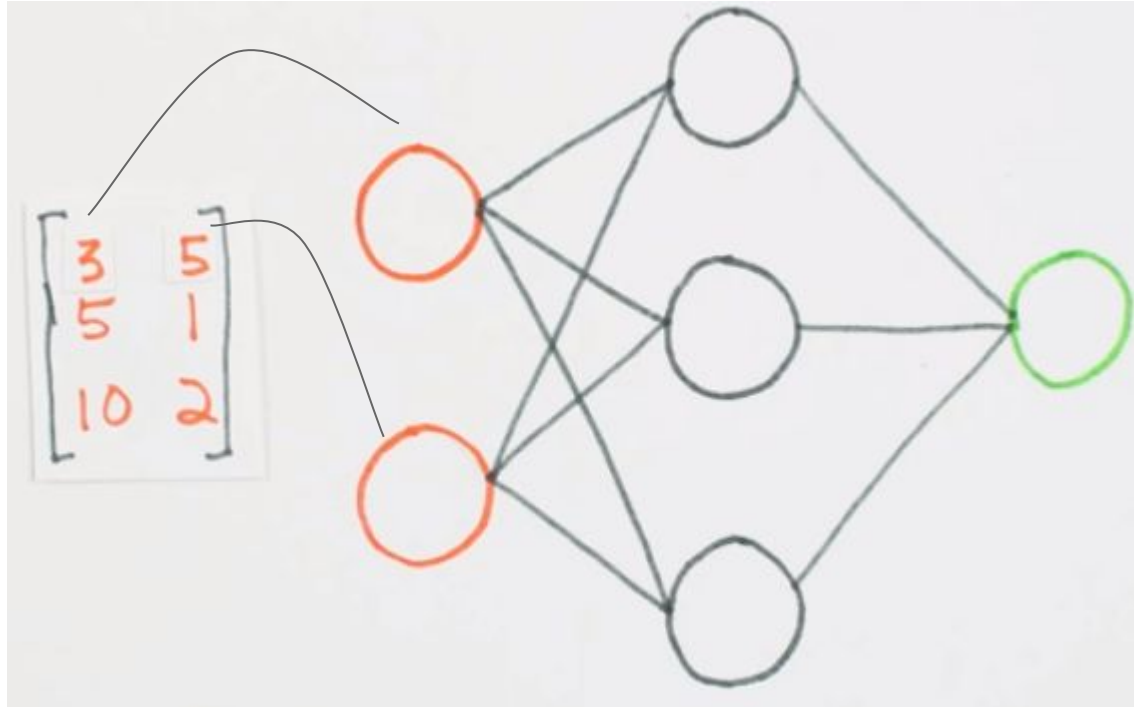


MULTI-NEURON NETWORKS :: ARCHITECTURE

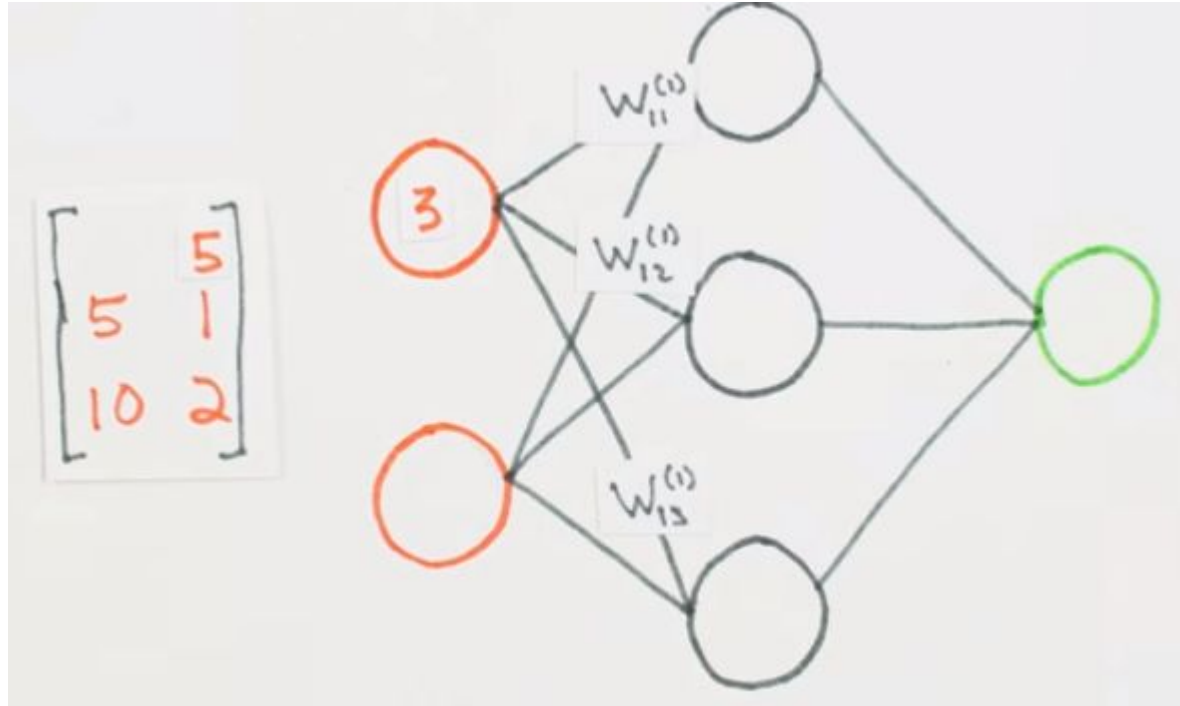
NEURON



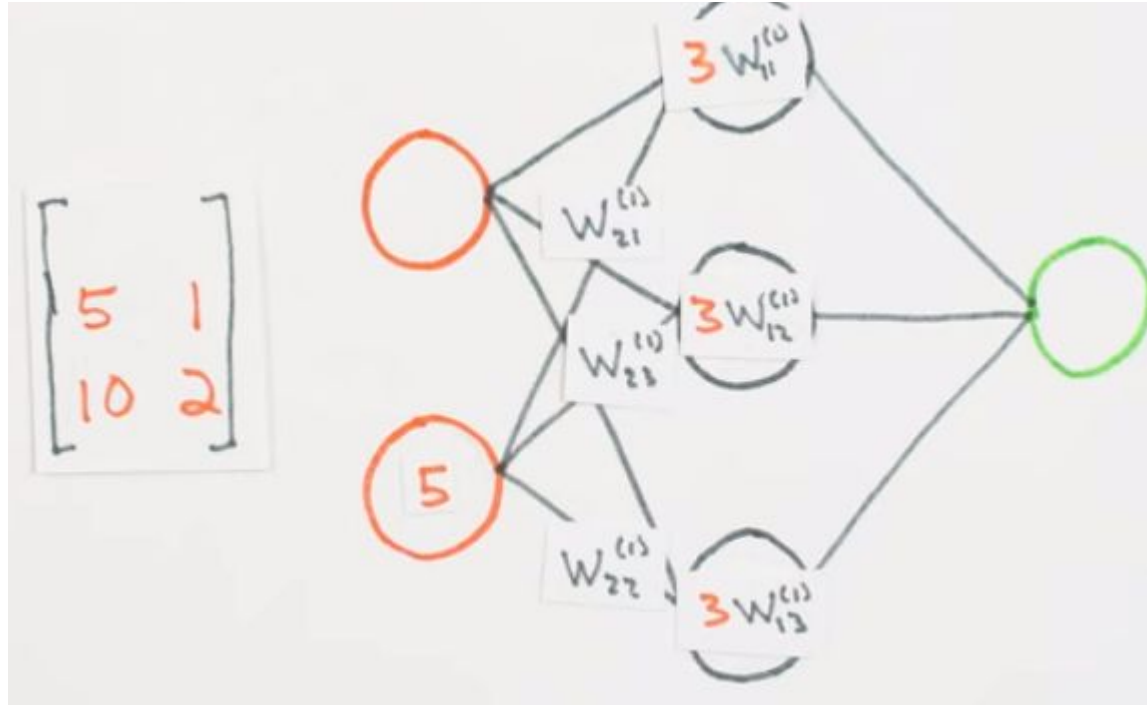
LEARNING STEP 1: FORWARD PROPAGATION



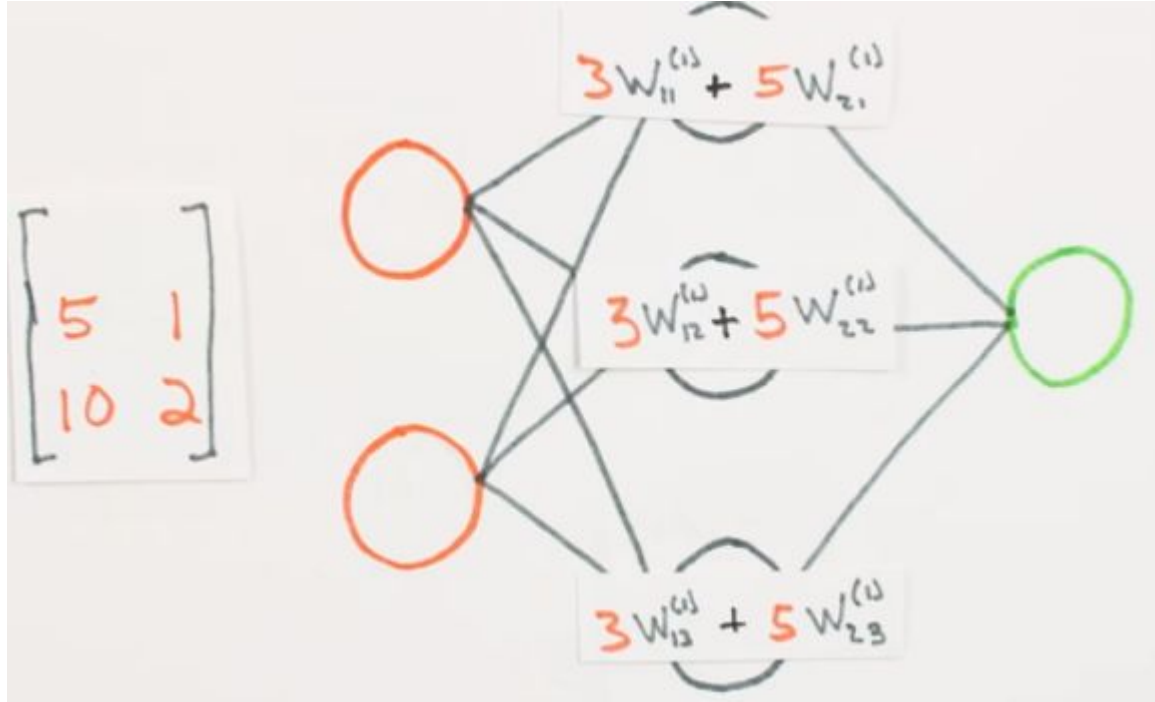
MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



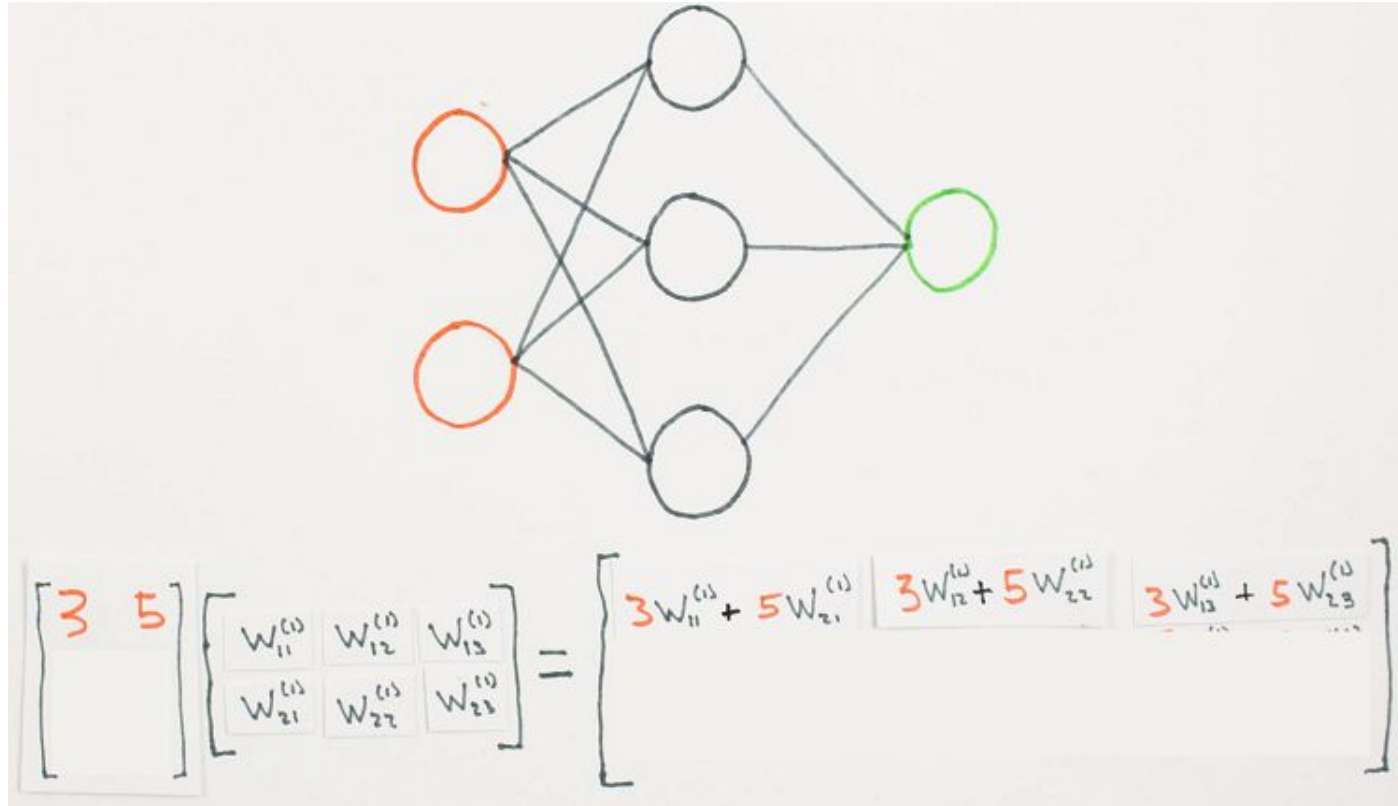
MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



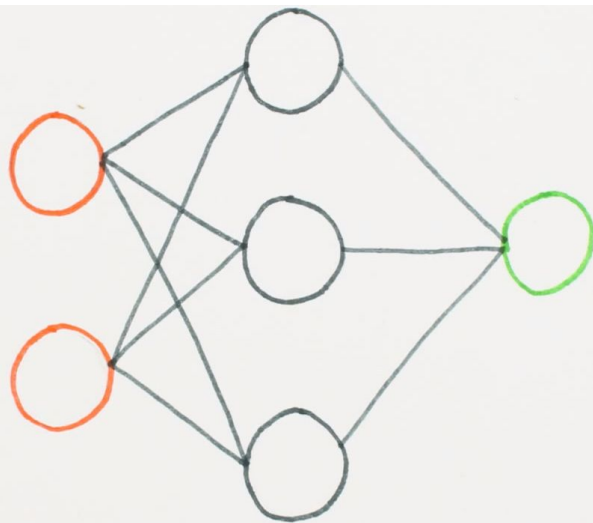
MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



MULTI-NEURON NETWORKS :: FORWARD PROPAGATION

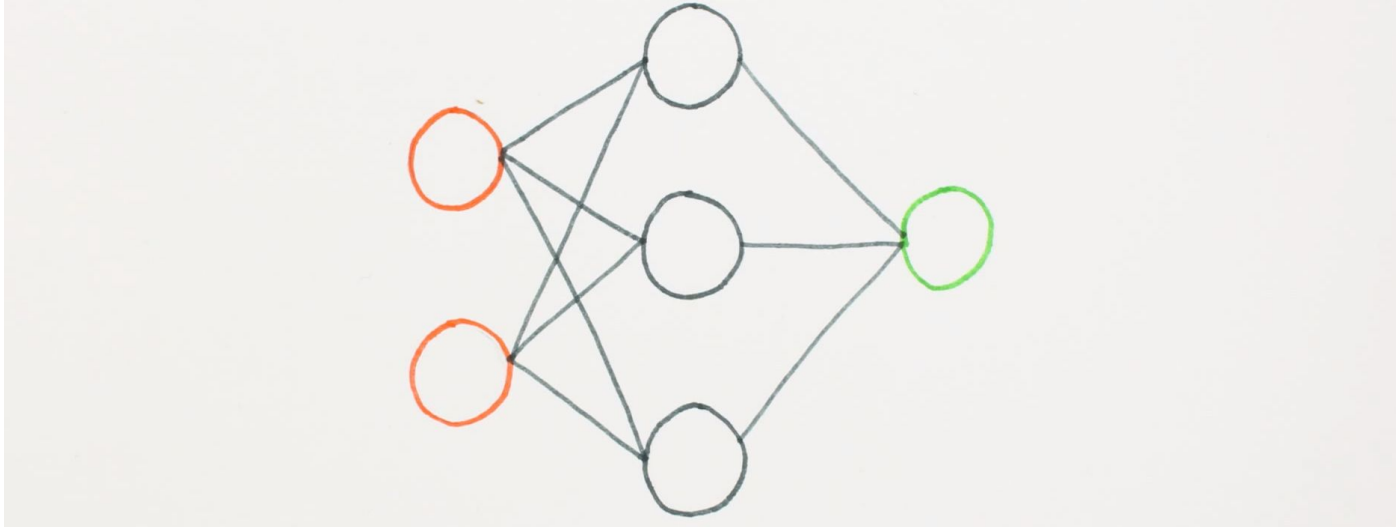


MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



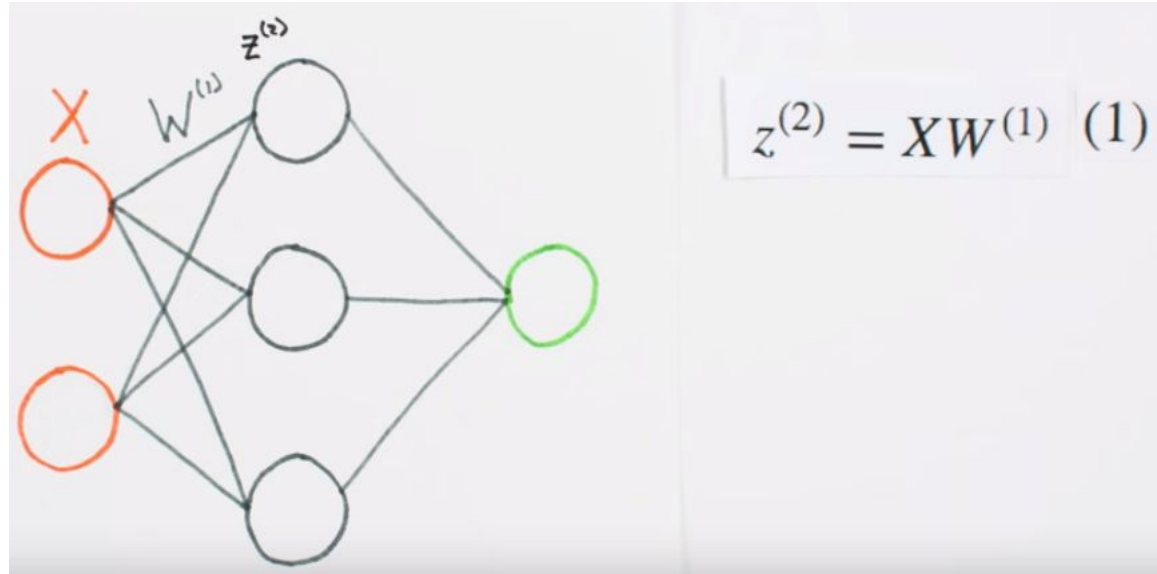
$$\begin{bmatrix} 3 & 5 \\ 5 & 1 \\ 10 & 2 \end{bmatrix} \begin{bmatrix} W_{11}^{(1)} & W_{12}^{(1)} & W_{13}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} & W_{23}^{(1)} \end{bmatrix} = \begin{bmatrix} 3W_{11}^{(1)} + 5W_{21}^{(1)} & 3W_{12}^{(1)} + 5W_{22}^{(1)} & 3W_{13}^{(1)} + 5W_{23}^{(1)} \\ 5W_{11}^{(1)} + 1W_{21}^{(1)} & 5W_{12}^{(1)} + 1W_{22}^{(1)} & 5W_{13}^{(1)} + 1W_{23}^{(1)} \\ 10W_{11}^{(1)} + 2W_{21}^{(1)} & 10W_{12}^{(1)} + 2W_{22}^{(1)} & 10W_{13}^{(1)} + 2W_{23}^{(1)} \end{bmatrix}$$

MULTI-NEURON NETWORKS :: FORWARD PROPAGATION

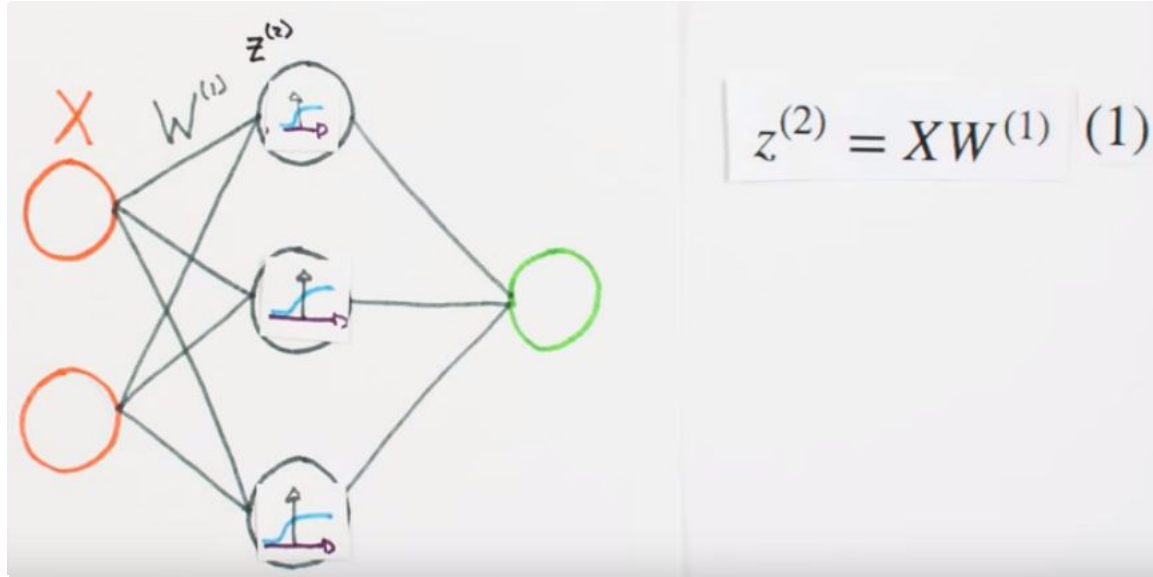


X	$W^{(1)}$	=	$Z^{(2)}$
---	-----------	---	-----------

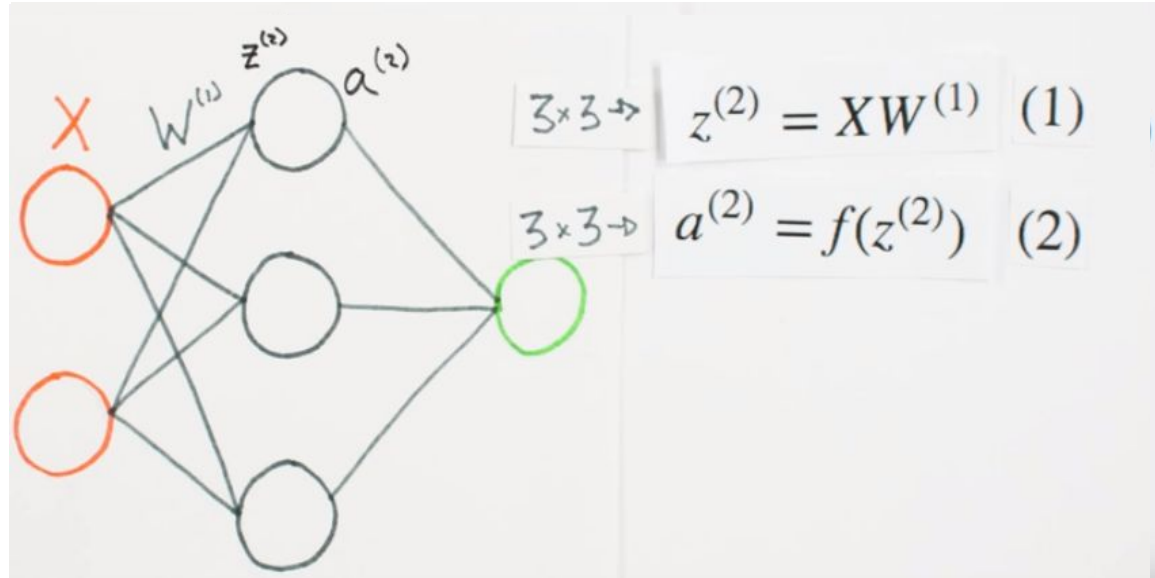
MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



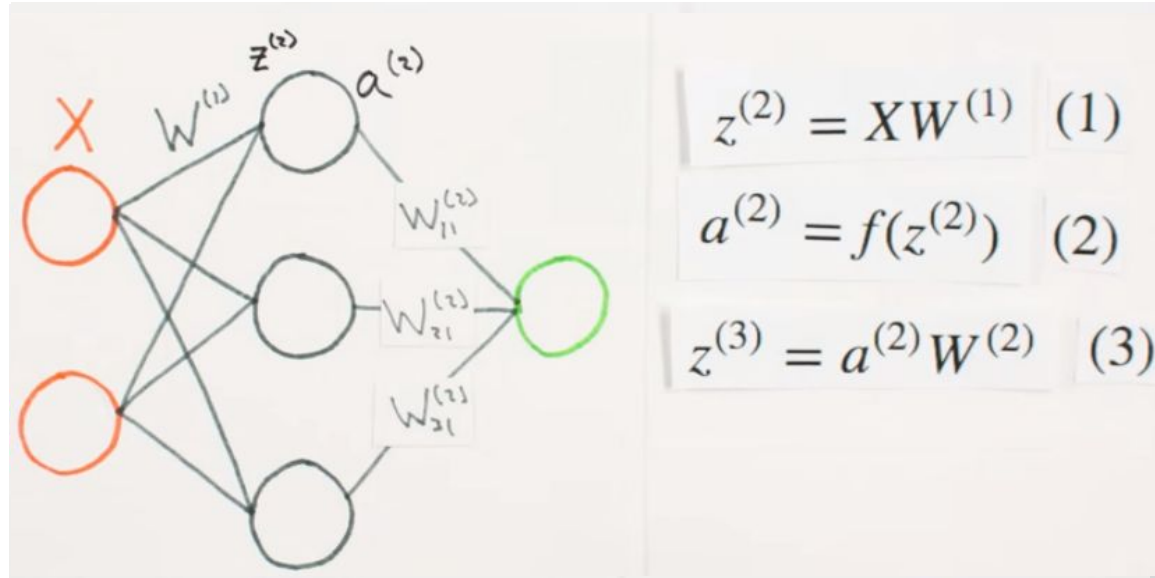
MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



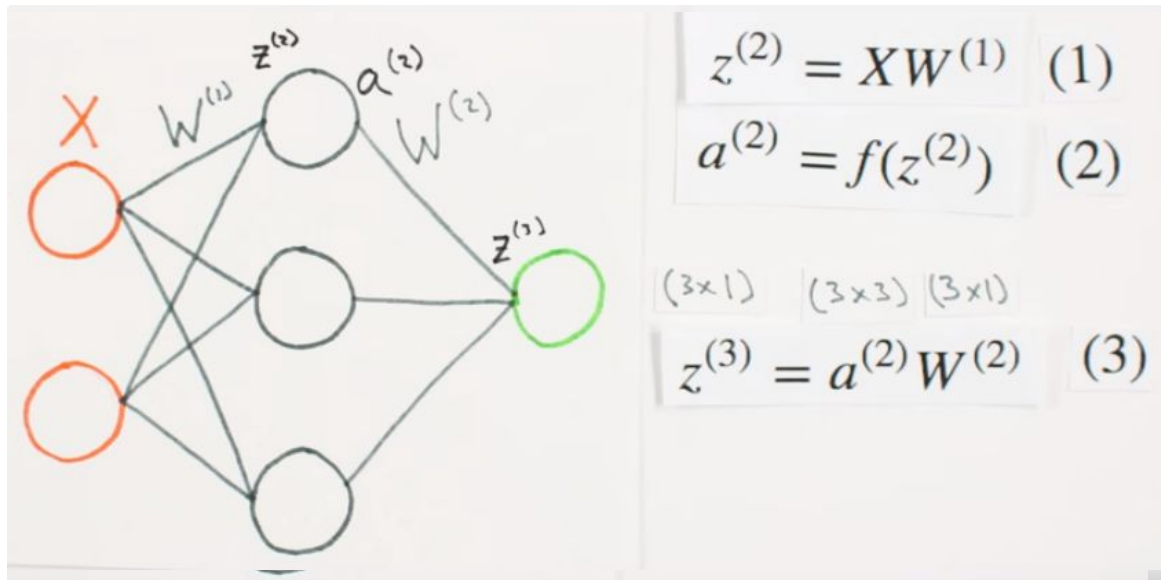
MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



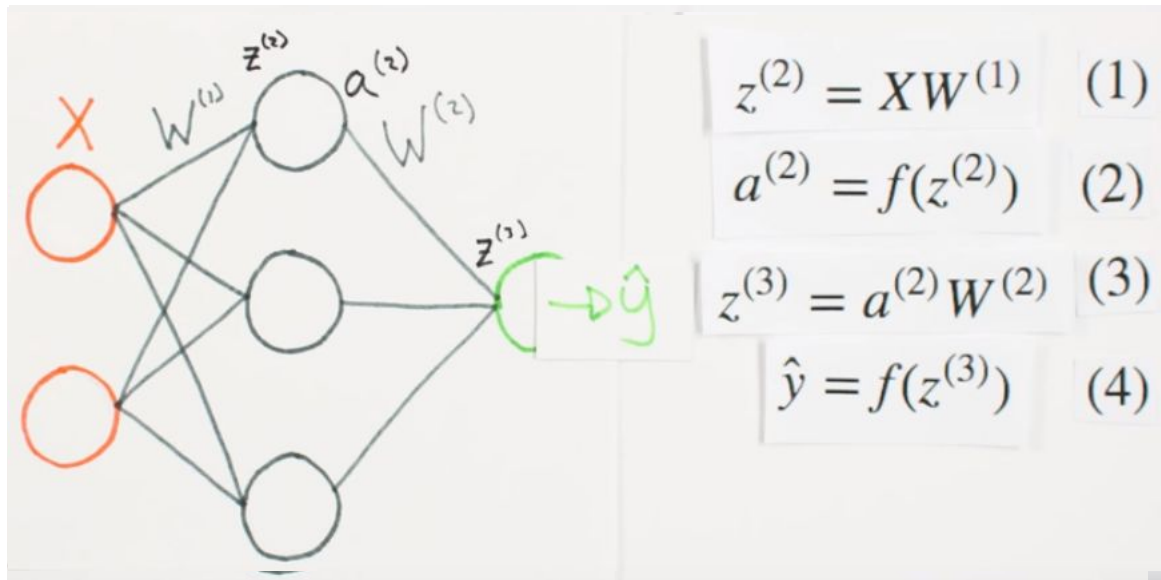
MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



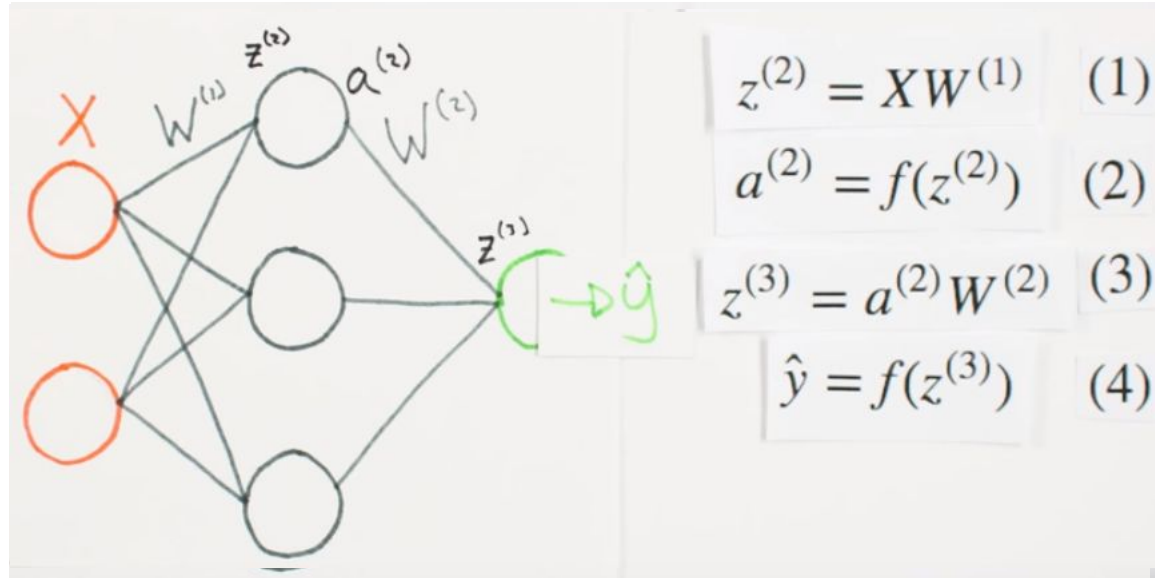
MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



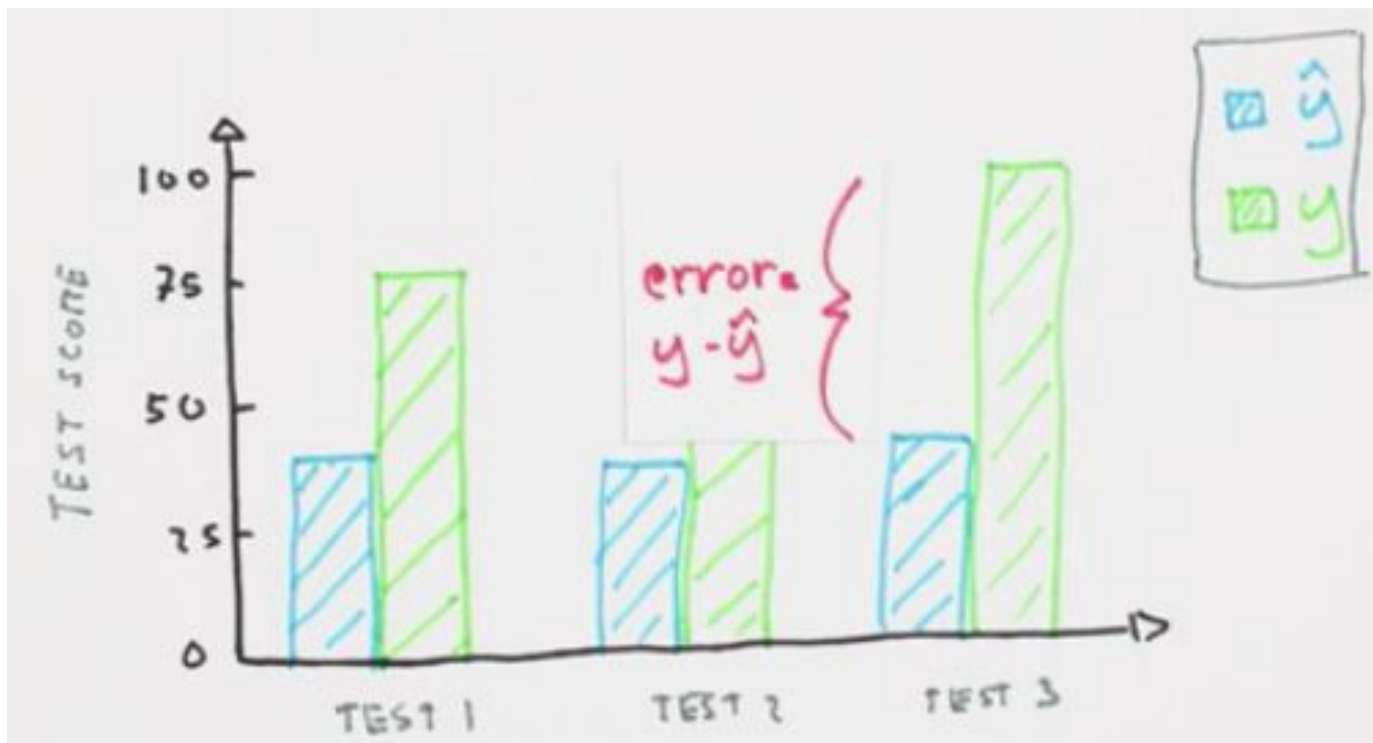
MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



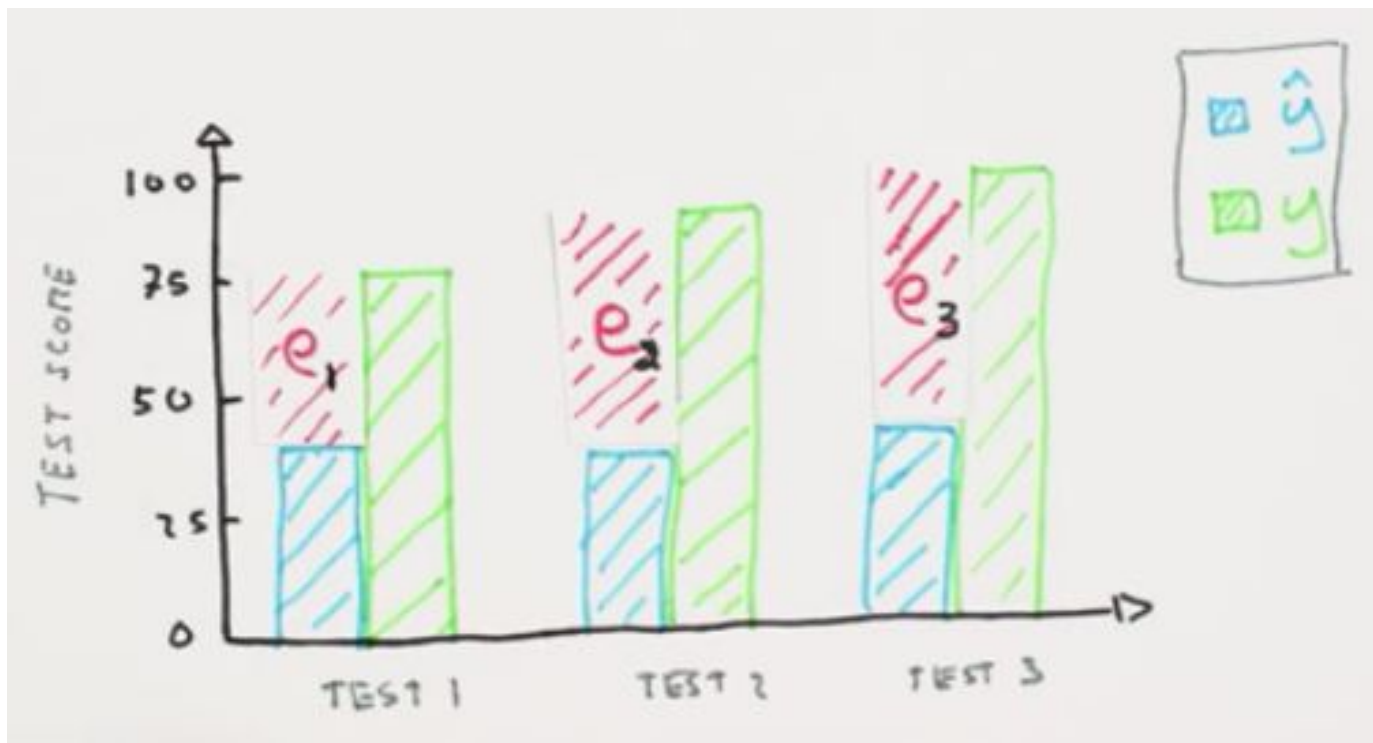
MULTI-NEURON NETWORKS :: FORWARD PROPAGATION



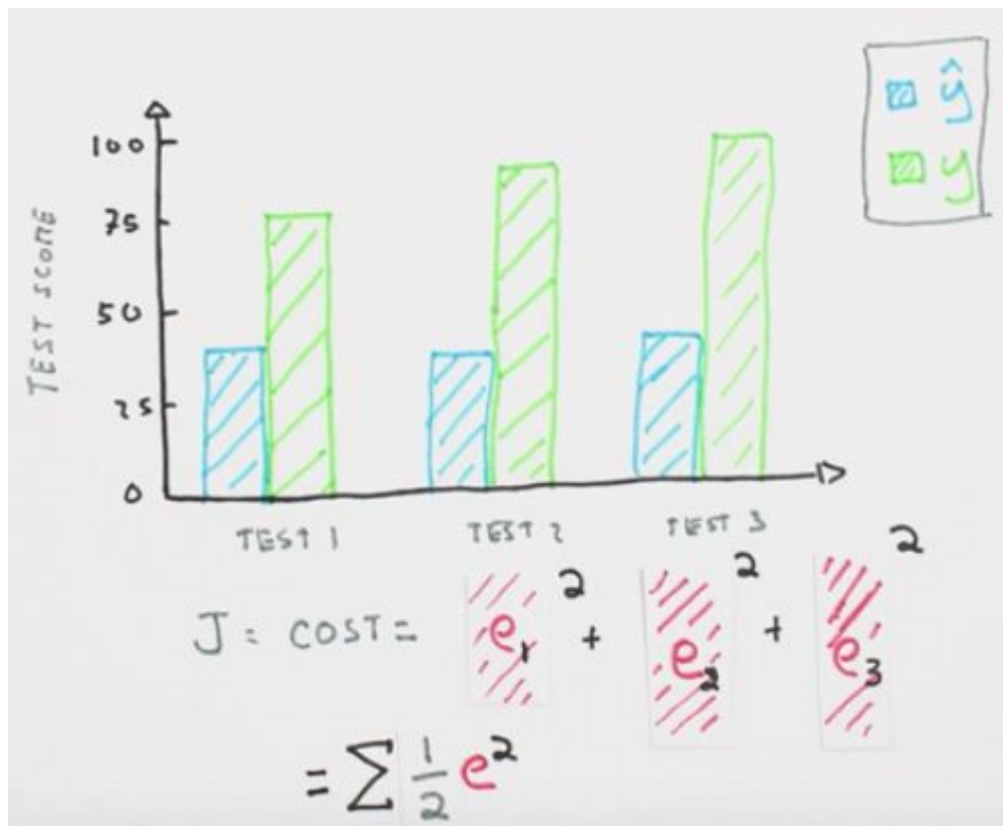
MULTI-NEURON NETWORKS :: GRADIENT DESCENT



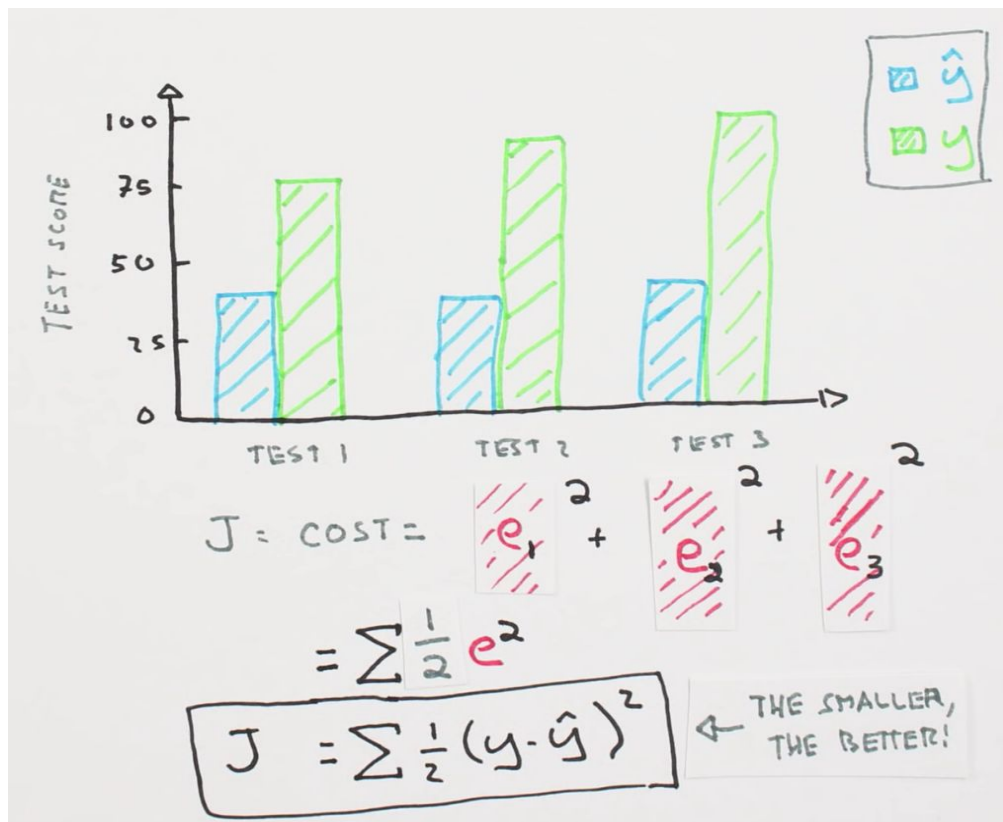
MULTI-NEURON NETWORKS :: GRADIENT DESCENT



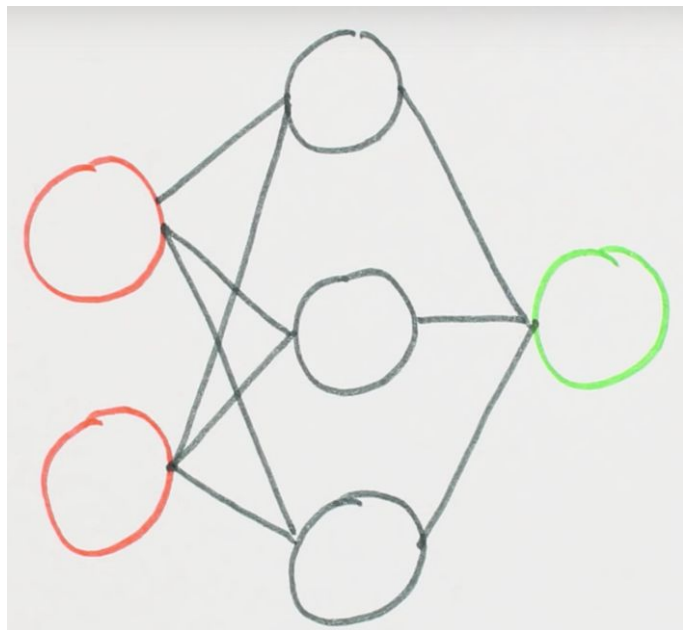
MULTI-NEURON NETWORKS :: GRADIENT DESCENT



MULTI-NEURON NETWORKS :: GRADIENT DESCENT

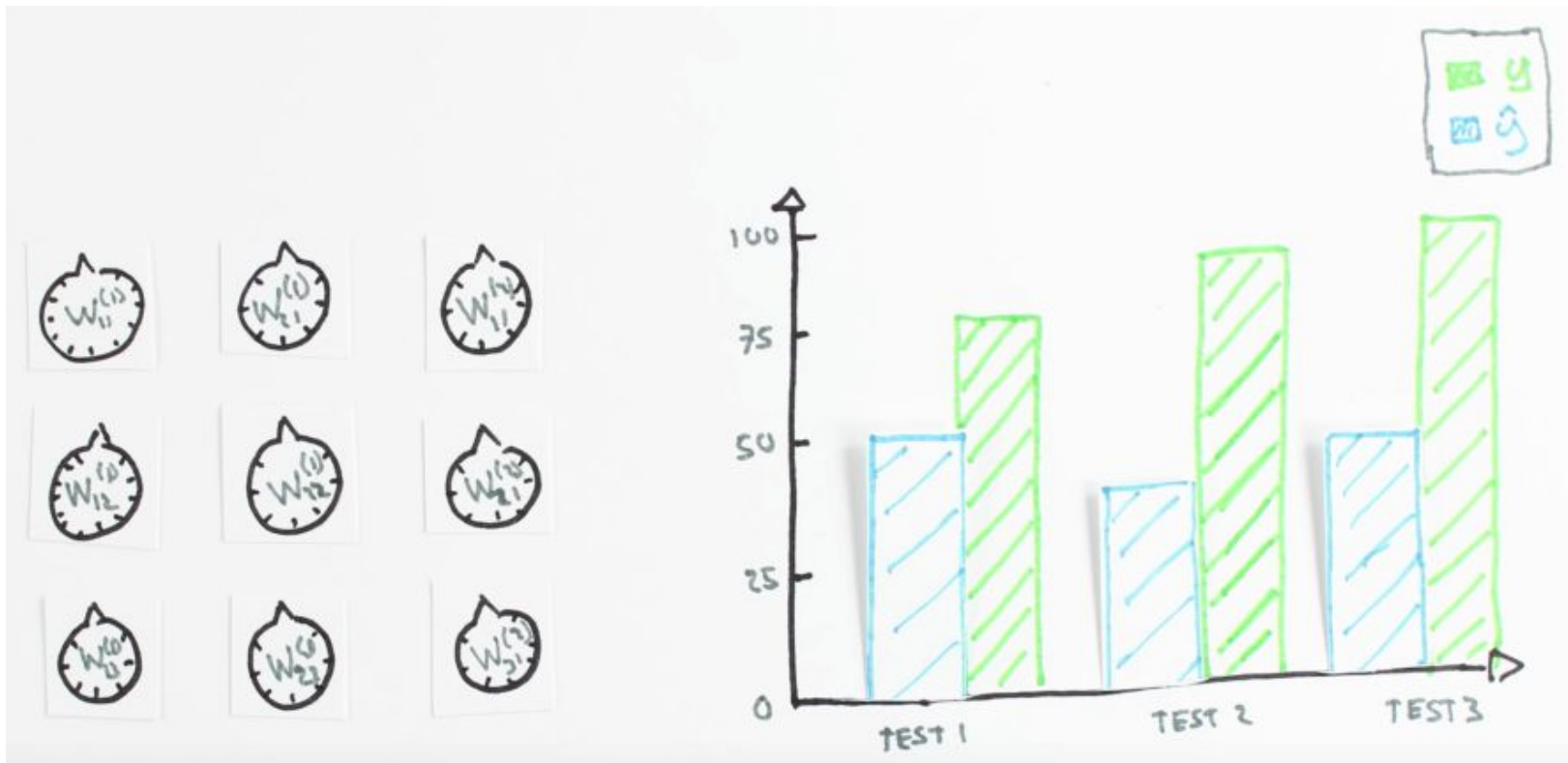


MULTI-NEURON NETWORKS :: GRADIENT DESCENT

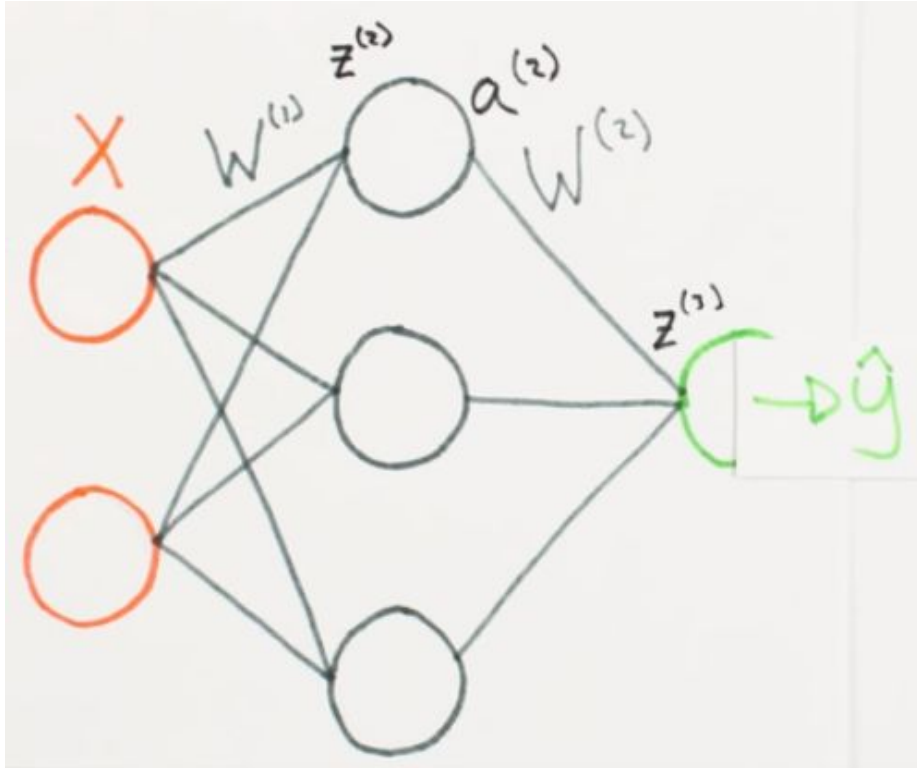


Training a Network
=
Minimizing a Cost
Function

MULTI-NEURON NETWORKS :: GRADIENT DESCENT



MULTI-NEURON NETWORKS :: GRADIENT DESCENT



$$z^{(2)} = XW^{(1)} \quad (1)$$

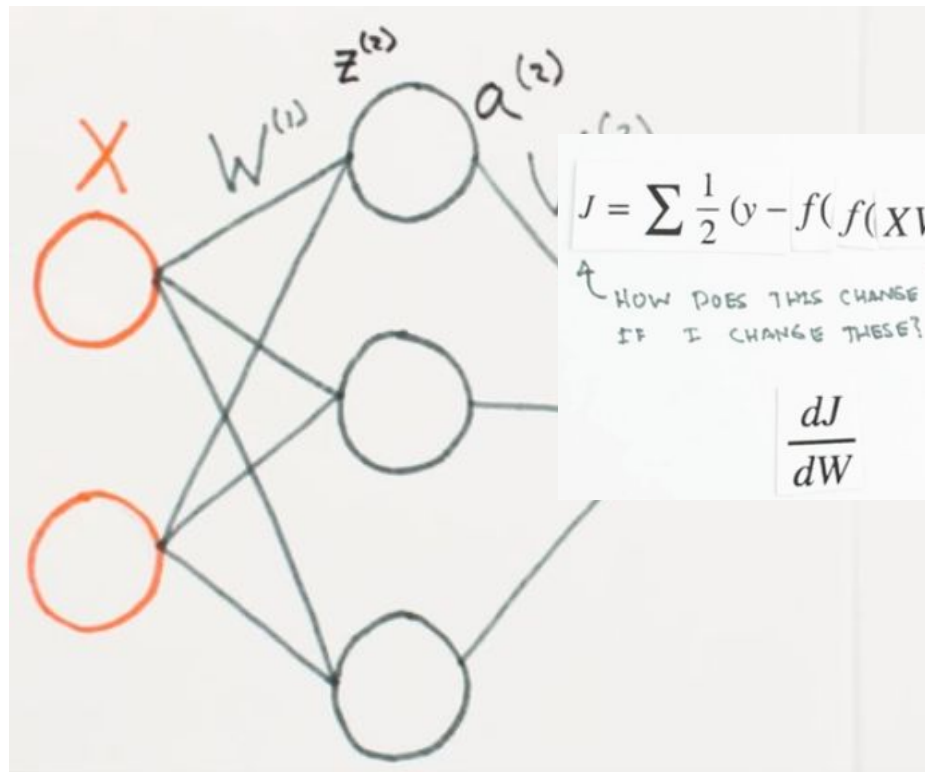
$$a^{(2)} = f(z^{(2)}) \quad (2)$$

$$z^{(3)} = a^{(2)}W^{(2)} \quad (3)$$

$$\hat{y} = f(z^{(3)}) \quad (4)$$

$$J = \sum \frac{1}{2} (y - \hat{y})^2 \quad (5)$$

MULTI-NEURON NETWORKS :: GRADIENT DESCENT



$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$

↑ HOW DOES THIS CHANGE
IF I CHANGE THESE?

$$\frac{dJ}{dW}$$

$$z^{(2)} = XW^{(1)} \quad (1)$$

$$a^{(2)} = f(z^{(2)}) \quad (2)$$

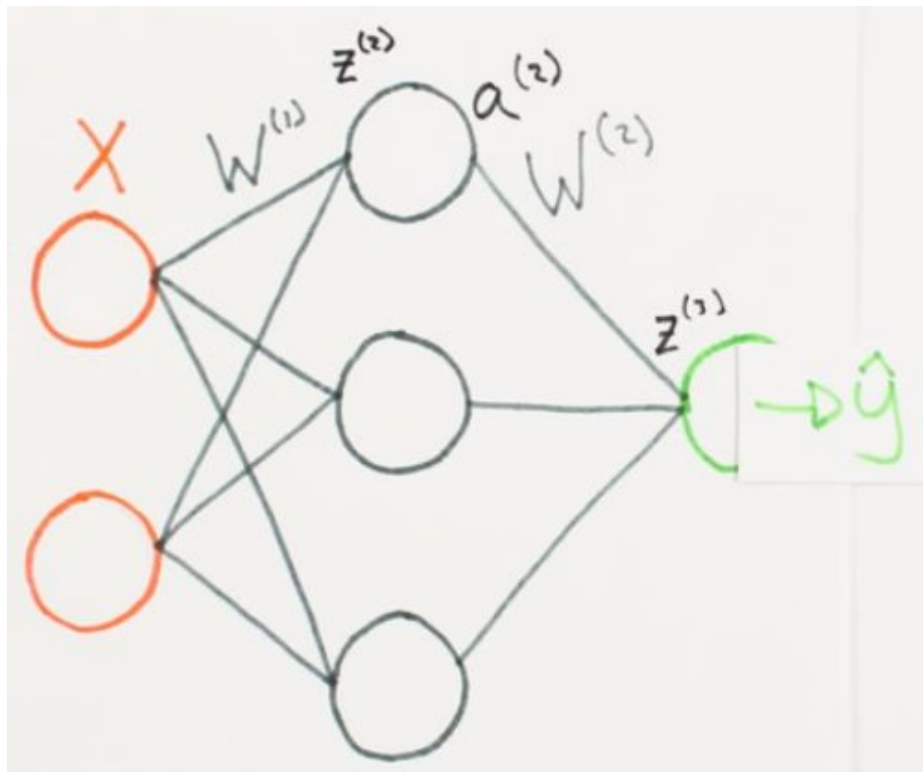
$$z^{(3)} = a^{(2)}W^{(2)} \quad (3)$$

$$\hat{y} = f(z^{(3)}) \quad (4)$$

$$J = \sum \frac{1}{2} (y - \hat{y})^2 \quad (5)$$

$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$

MULTI-NEURON NETWORKS :: GRADIENT DESCENT

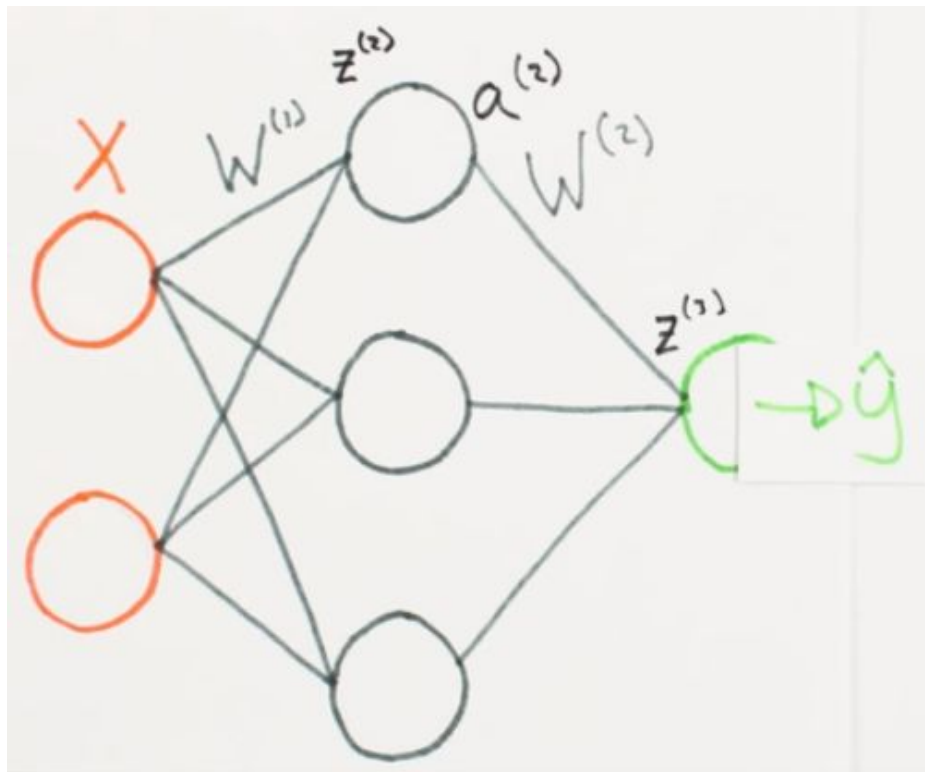


$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$

↑ HOW DOES THIS CHANGE
IF I CHANGE THESE? ↑

$$\frac{dJ}{dW}$$

MULTI-NEURON NETWORKS :: GRADIENT DESCENT

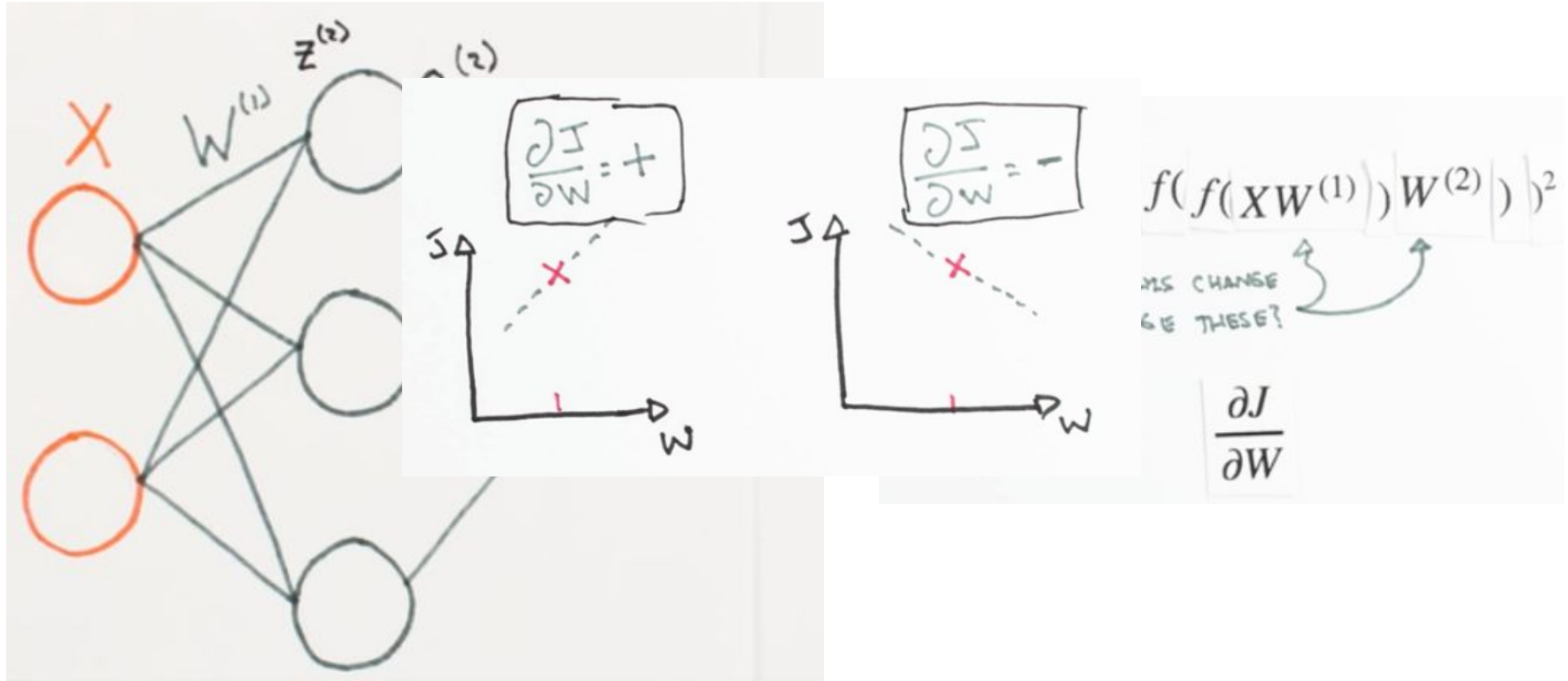


$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$

↑ HOW DOES THIS CHANGE
IF I CHANGE THESE? ↑

$$\frac{\partial J}{\partial W}$$

MULTI-NEURON NETWORKS :: GRADIENT DESCENT

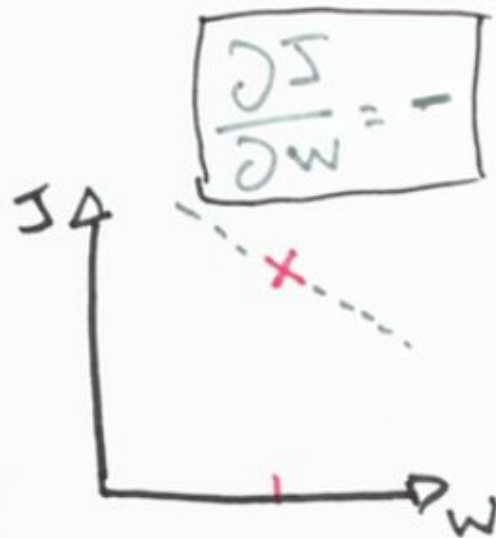
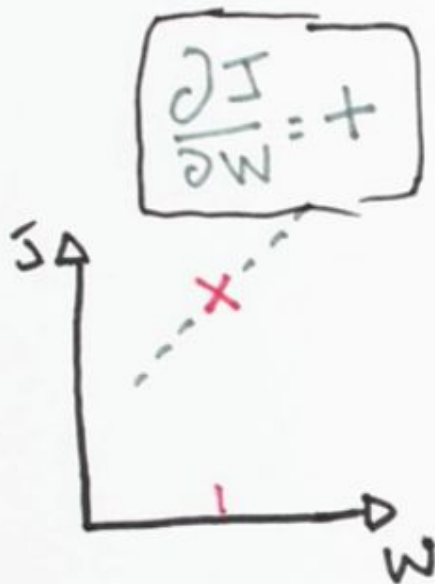


MULTI-NEURON NETWORKS :: GRADIENT DESCENT

$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$

↑ HOW DOES THIS CHANGE
IF I CHANGE THESE?

$$\frac{\partial J}{\partial W}$$

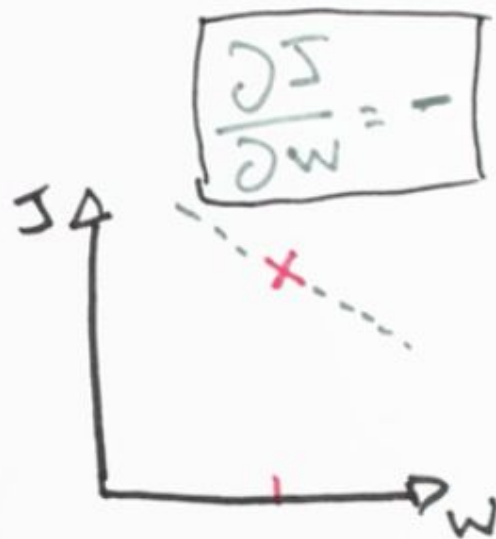
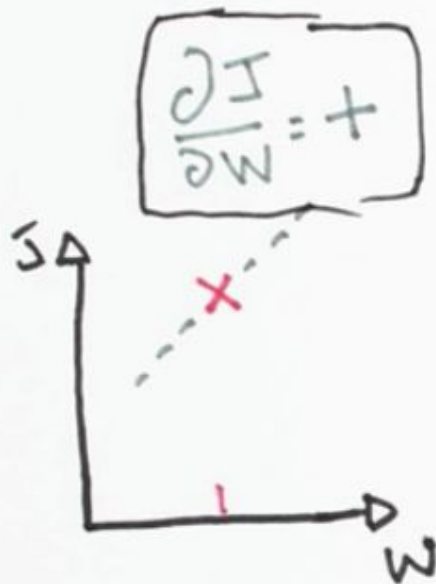


MULTI-NEURON NETWORKS :: GRADIENT DESCENT

$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$

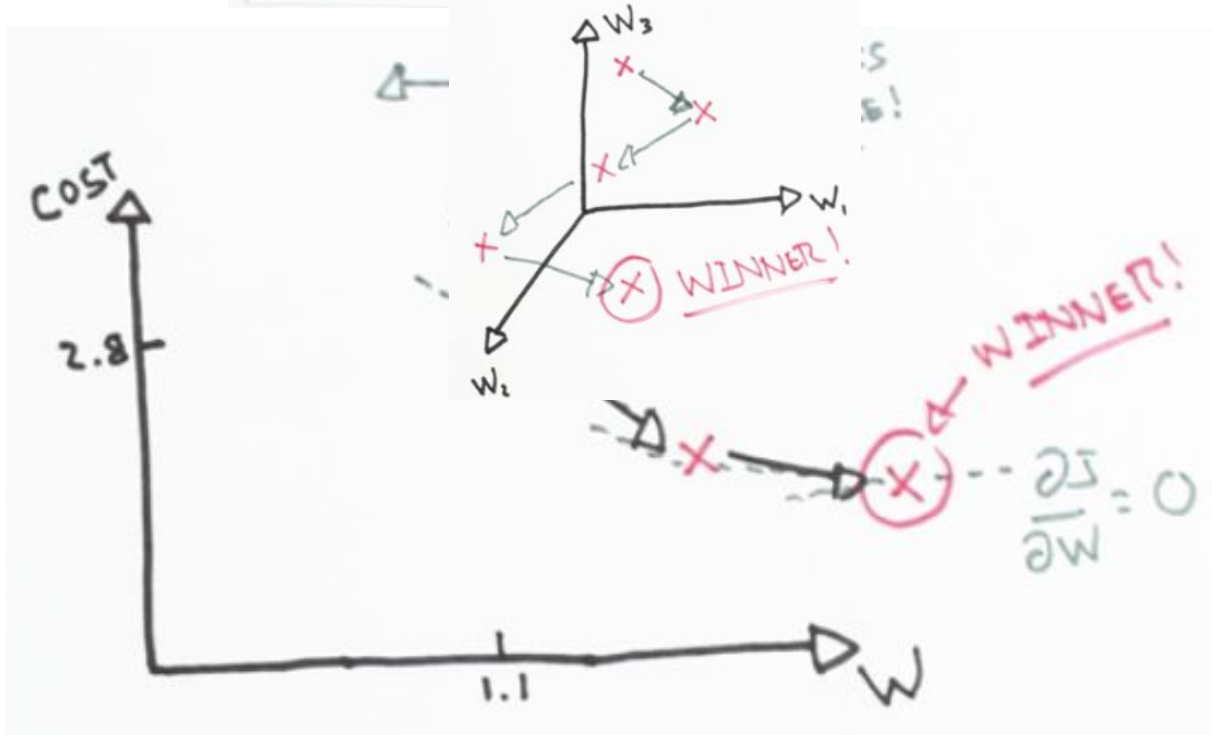
↑ HOW DOES THIS CHANGE
IF I CHANGE THESE?

$$\frac{\partial J}{\partial W}$$



MULTI-NEURON NETWORKS :: GRADIENT DESCENT

$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$

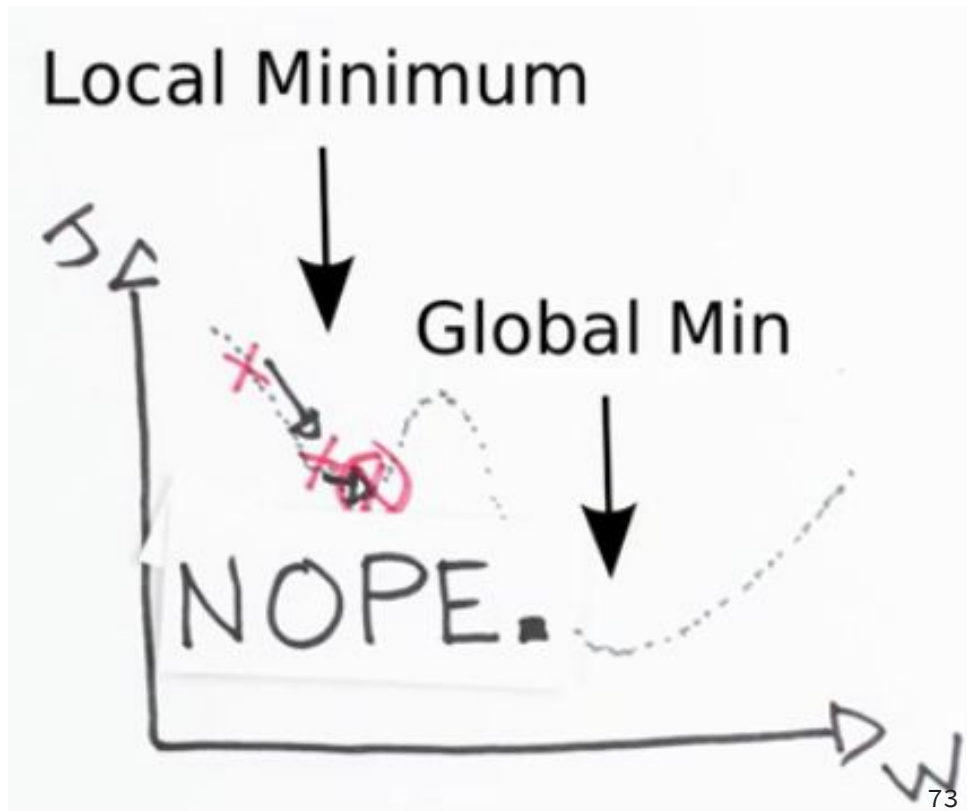


MULTI-NEURON NETWORKS :: GRADIENT DESCENT

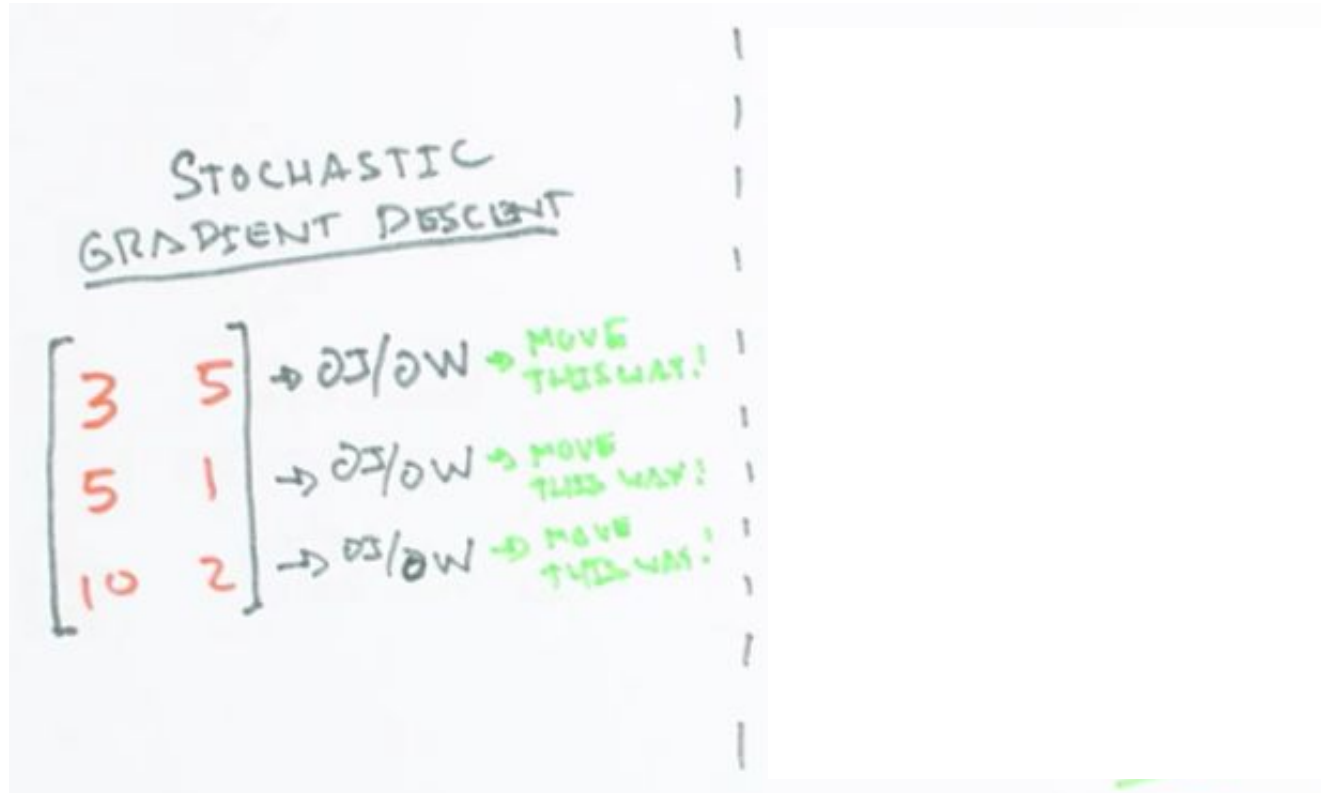


MULTI-NEURON NETWORKS :: GRADIENT DESCENT

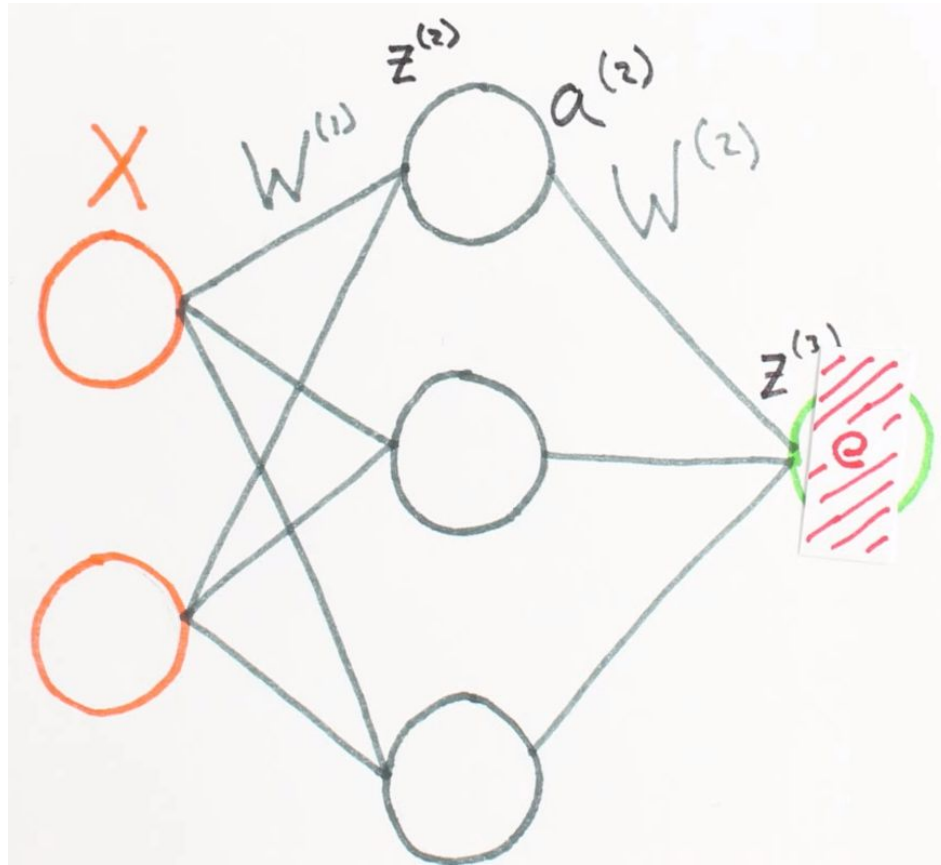
$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2$$



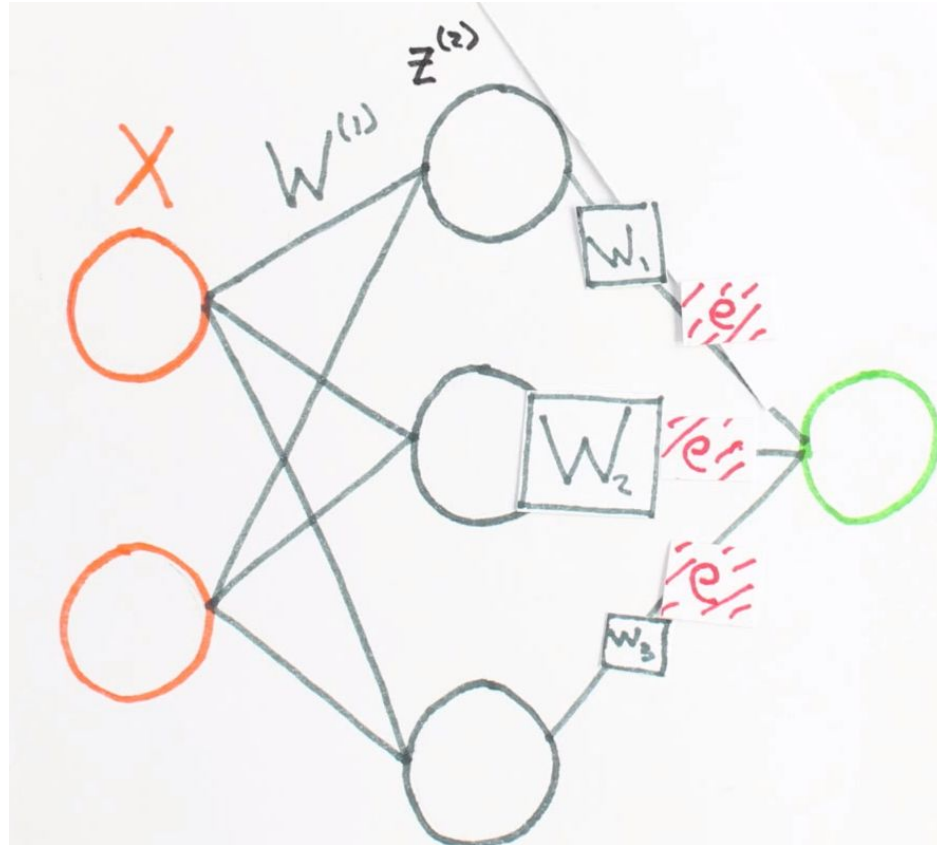
MULTI-NEURON NETWORKS :: GRADIENT DESCENT



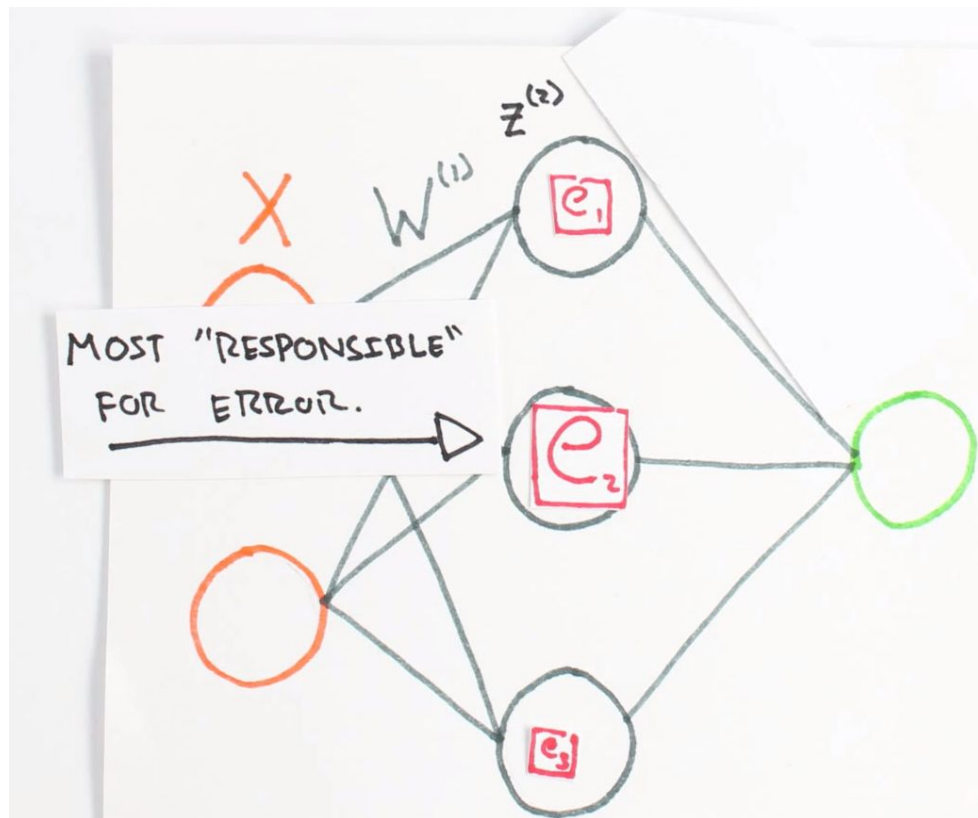
MULTI-NEURON NETWORKS :: GRADIENT DESCENT



MULTI-NEURON NETWORKS :: GRADIENT DESCENT



MULTI-NEURON NETWORKS :: GRADIENT DESCENT



MULTI-NEURON NETWORKS :: TRAINING

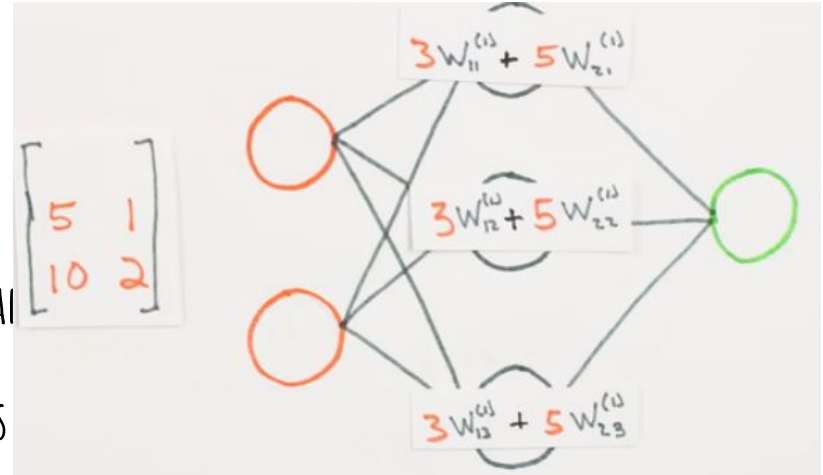
INITIALIZE NETWORK WITH RANDOM WEIGHTS

WHILE [NOT CONVERGED]

DO FORWARD PROP

DO BACKPROP AND DETERMINE CHA

UPDATE ALL WEIGHTS IN ALL LAYERS



MULTI-NEURON NETWORKS :: TRAINING

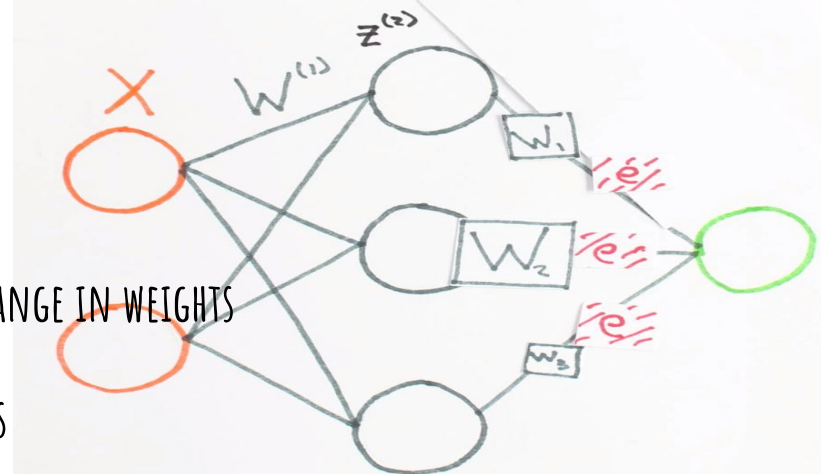
INITIALIZE NETWORK WITH RANDOM WEIGHTS

WHILE [NOT CONVERGED]

DO FORWARD PROP

DO BACKPROP AND DETERMINE CHANGE IN WEIGHTS

UPDATE ALL WEIGHTS IN ALL LAYERS



MULTI-NEURON NETWORKS :: TRAINING

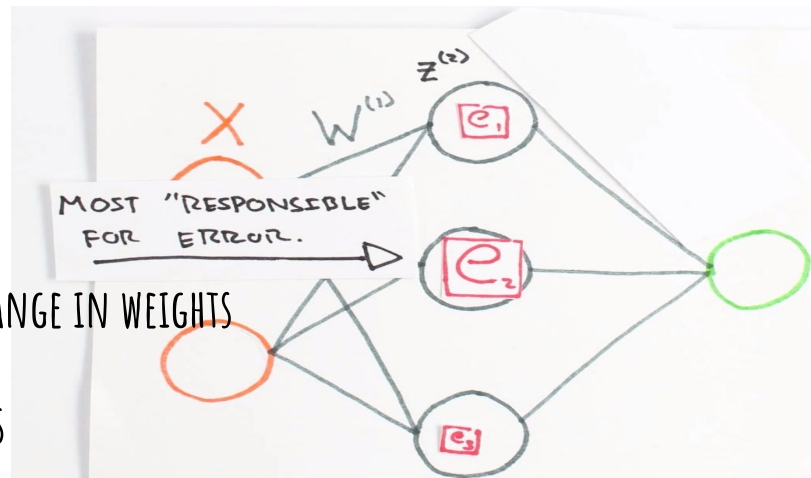
INITIALIZE NETWORK WITH RANDOM WEIGHTS

WHILE [NOT CONVERGED]

DO FORWARD PROP

DO BACKPROP AND DETERMINE CHANGE IN WEIGHTS

UPDATE ALL WEIGHTS IN ALL LAYERS



MULTI-NEURON NETWORKS :: TRAINING

INITIALIZE NETWORK WITH RANDOM WEIGHTS

WHILE [NOT CONVERGED]

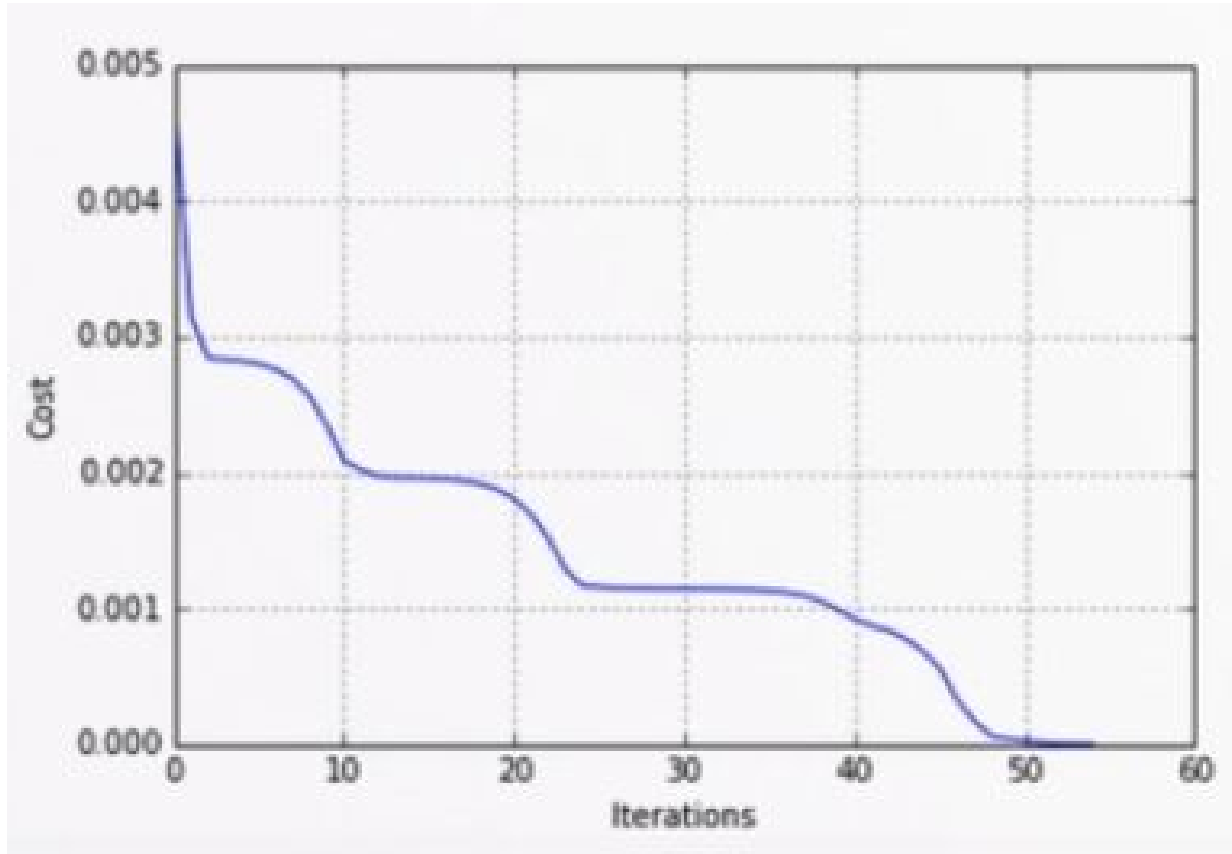
DO FORWARD PROP

DO BACKPROP AND DETERMINE CHANGE IN WEIGHTS

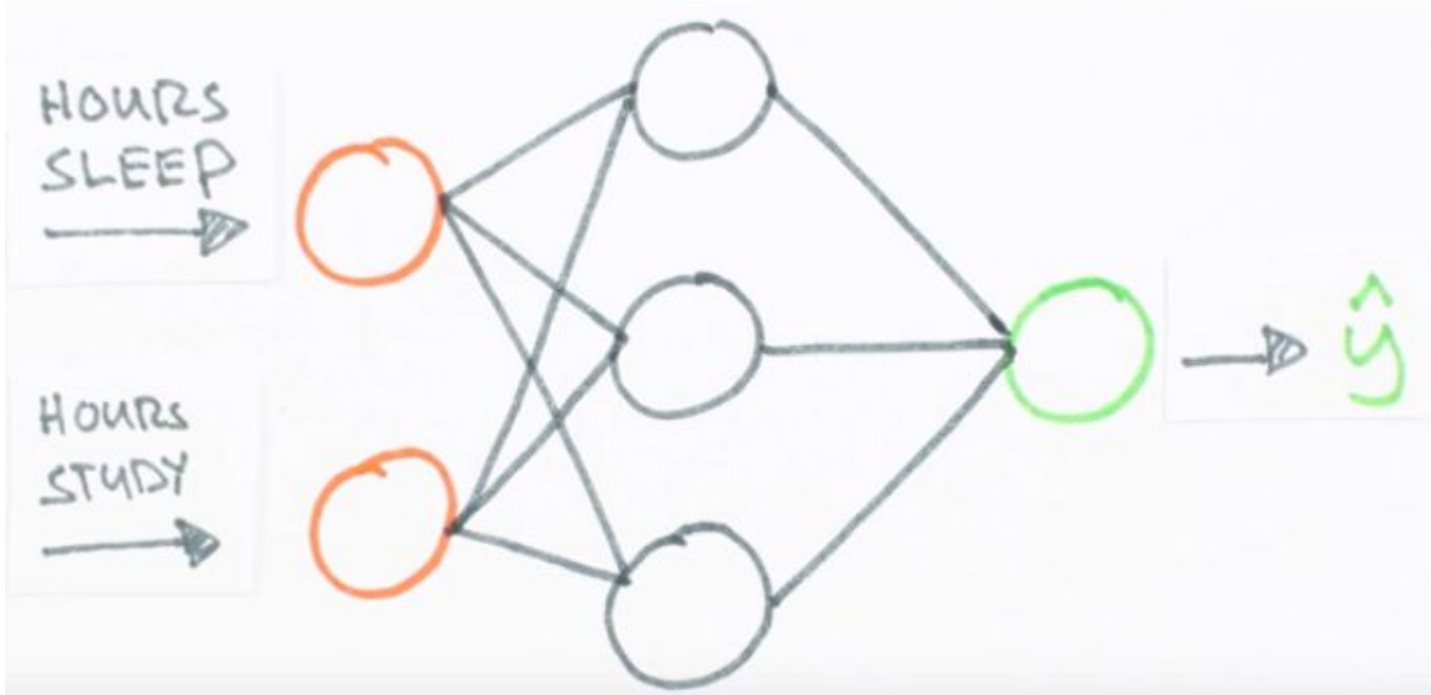
UPDATE ALL WEIGHTS IN ALL LAYERS

One
Iteration

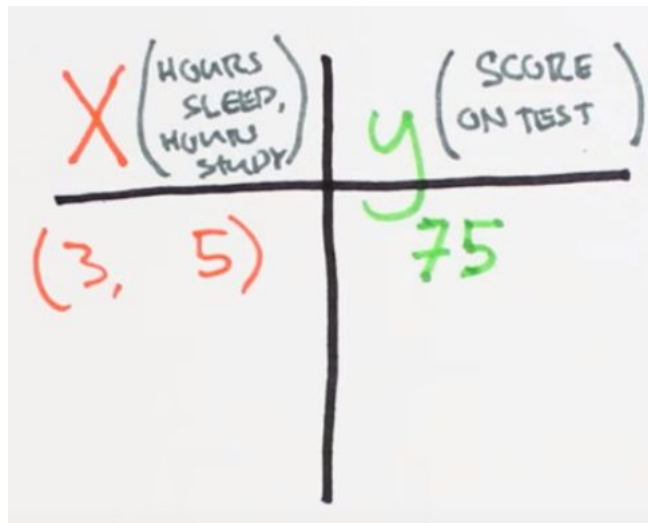
MULTI-NEURON NETWORKS :: TRAINING



MULTI-NEURON NETWORKS :: TESTING

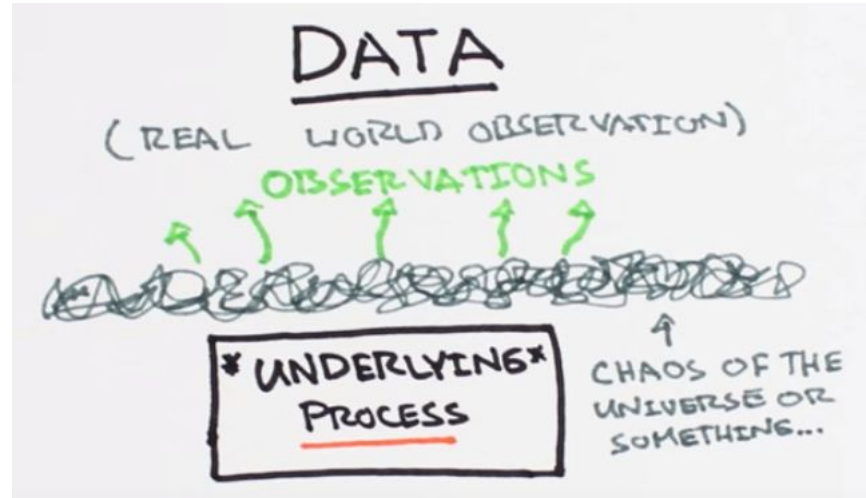


ALL IZZ WELL ???

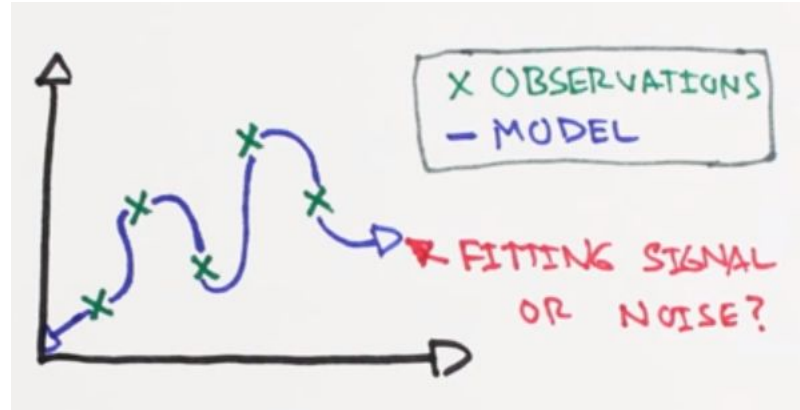


Demo ???

GENERAL SITUATION



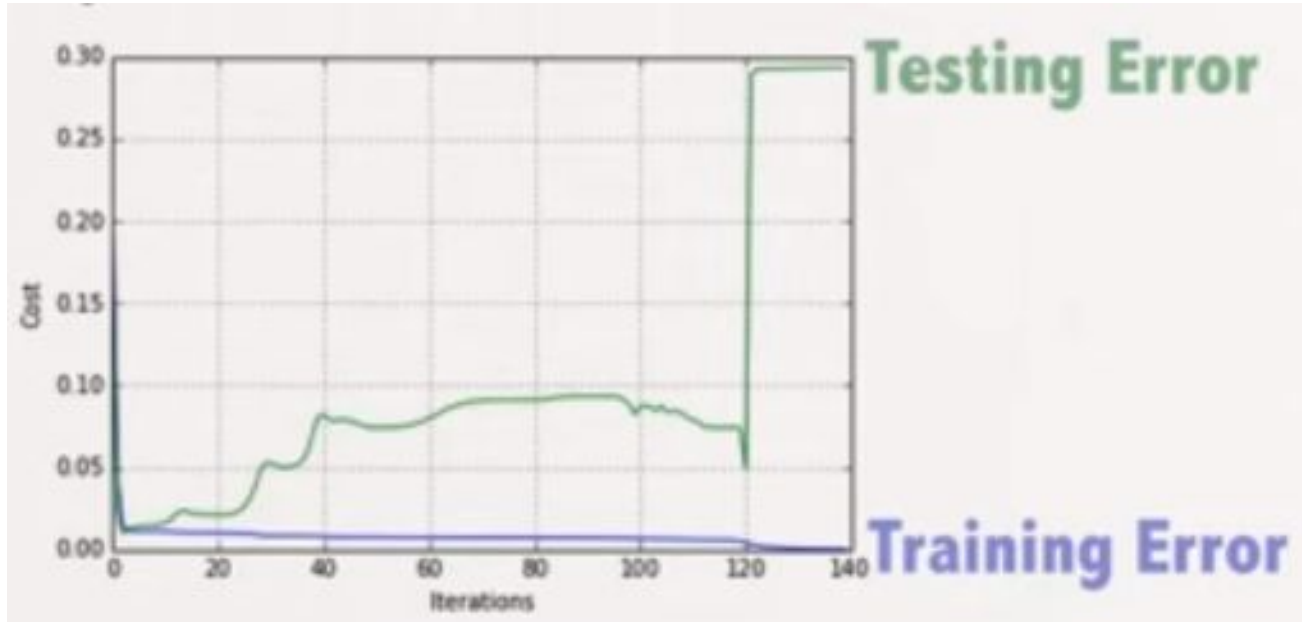
$$\text{OBSERVATION} = \text{SIGNAL} + \text{NOISE}$$



HOW TO AVOID OVERFITTING ?



HOW TO AVOID OVERFITTING ?

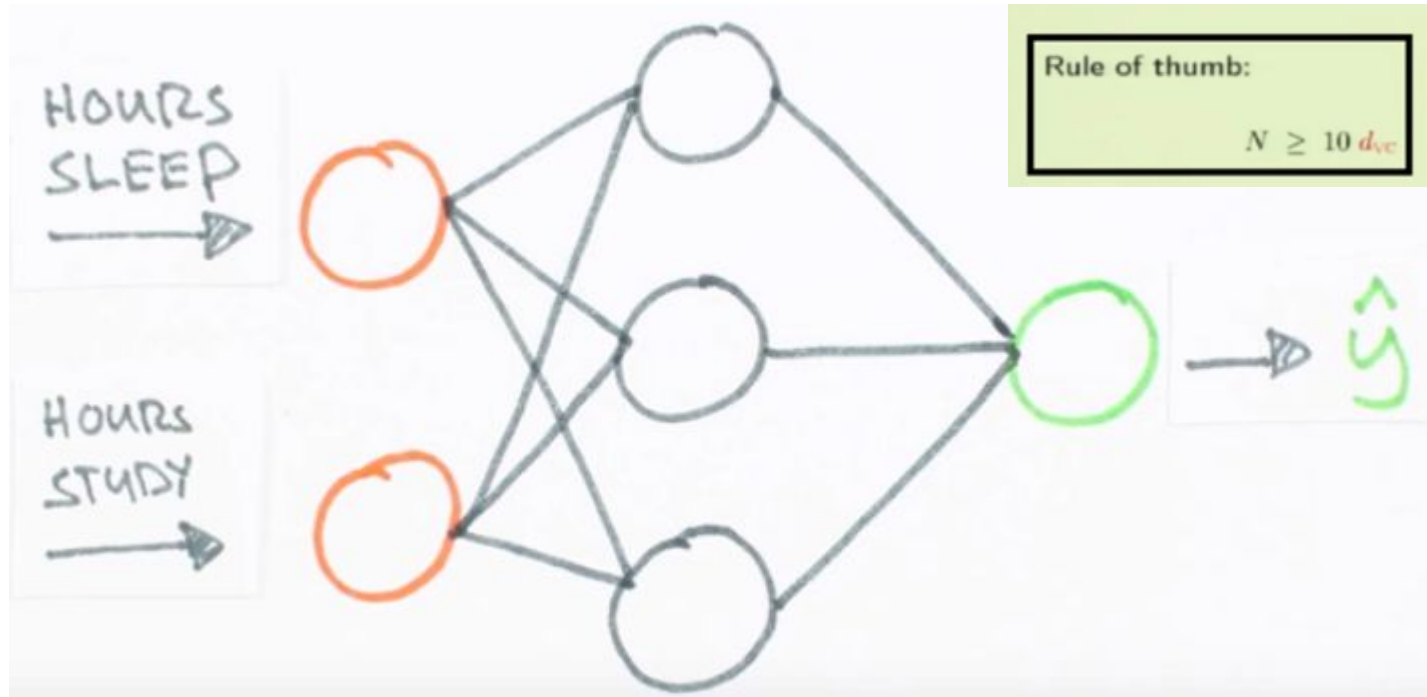


HOW MUCH DATA DO WE NEED ?

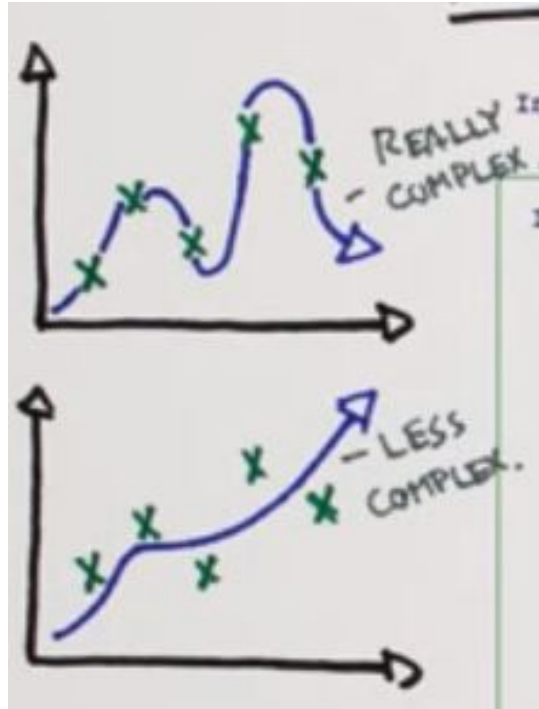
Rule of thumb:

$$N \geq 10 d_{vc}$$

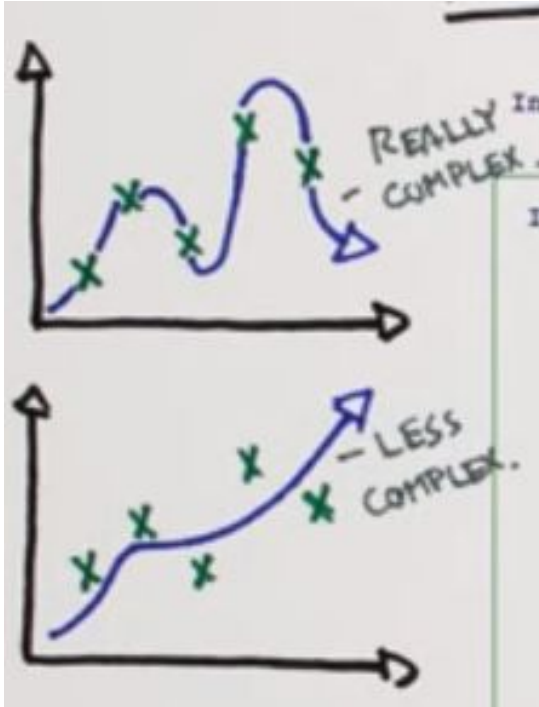
HOW MUCH DATA DO WE NEED ?



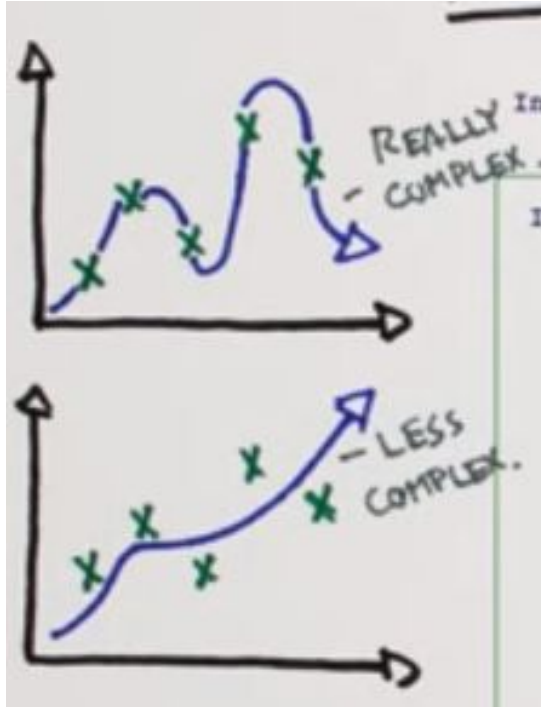
AVOIDING OVERFITTING : TECHNIQUE #2 : REGULARIZATION



AVOIDING OVERFITTING : TECHNIQUE #2 : REGULARIZATION

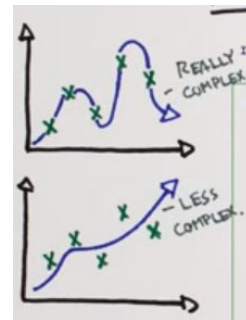
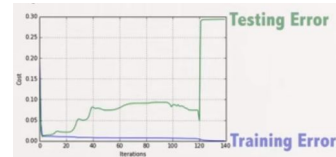
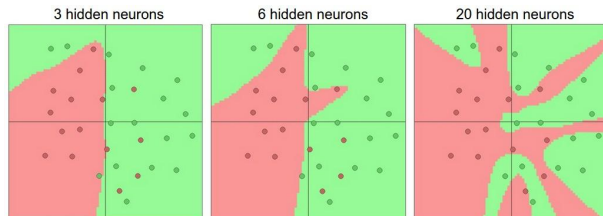
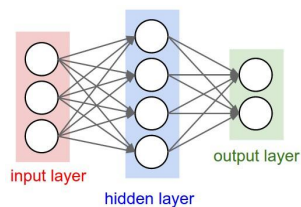
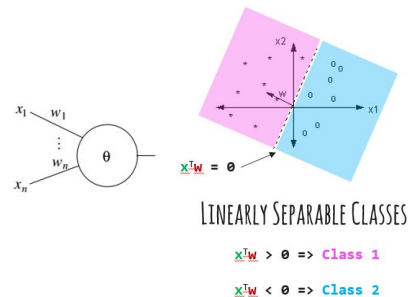
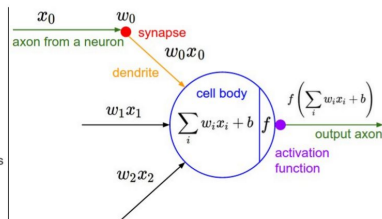
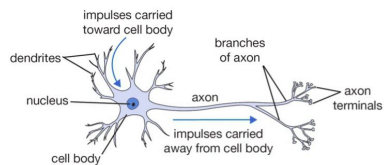


AVOIDING OVERFITTING : TECHNIQUE #2 : REGULARIZATION



$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)})W^{(2)}))^2 + ||W||^2$$

NEURAL NETWORKS



**3 HOURS OF NEURAL
NETWORKS !**

MIND = BLOWN

REFERENCES

- Neural network diagrams
<https://github.com/battlesnake/neural>
- ‘Neural Networks Demystified’ – YouTube channel
- Andrej Karpathy + Stanford + Neural Networks
- <http://neuralnetworksanddeeplearning.com/>