

Our example will be to prove Moore's Law.

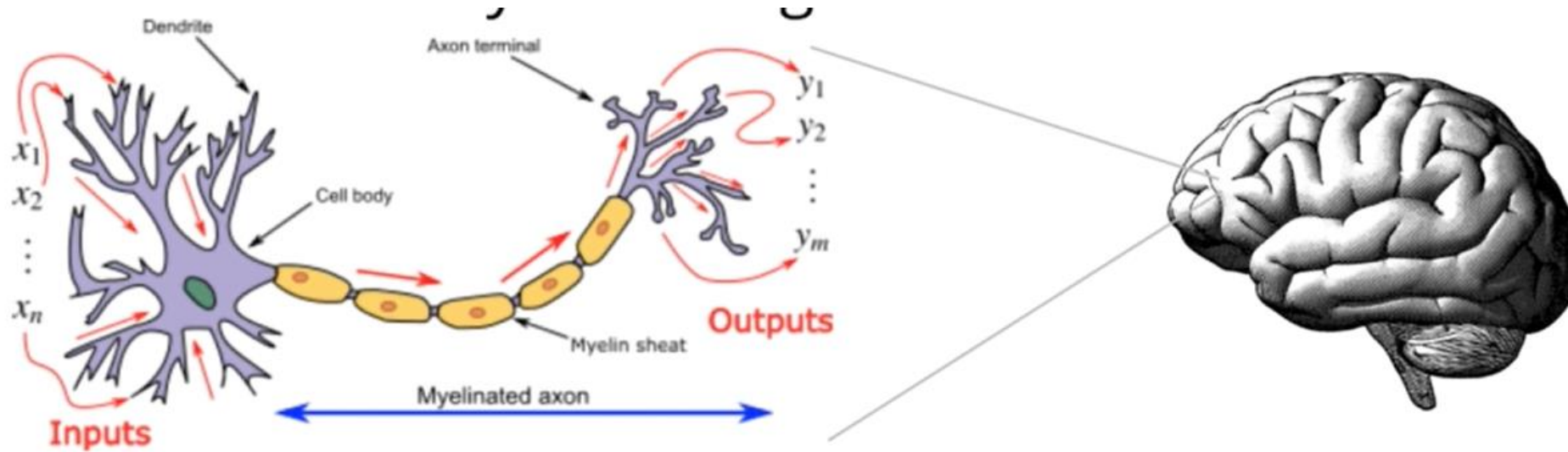
## • Refer NB “2 Linear\_Regression.ipynb”

Q.) What is Moore law?

Soln: States that the # Transistors / sq. inch on Integrated Circuits  
doubles approx. every 2 yrs.  $\Rightarrow$  Exponential Growth  
 $\downarrow$   
Standardization (Log)

# Artificial Neural Network (ANN)

# Where do ANN come from?



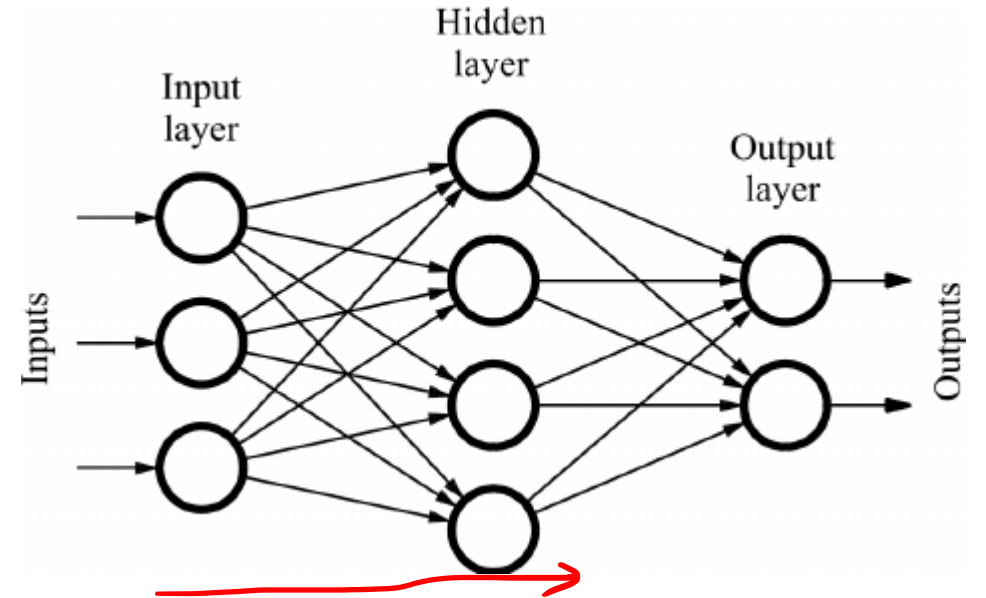
# The obvious question

- We know that brain is made up of neurons, and we also know how neurons work (and thus can simulate them)
- It thus makes sense to ask: “Can we build brain?”
- If we connect a bunch of neurons together, will intelligence suddenly emerge?
- If so, it would be an Artificial Intelligence.

# Feed Forward Neural Network

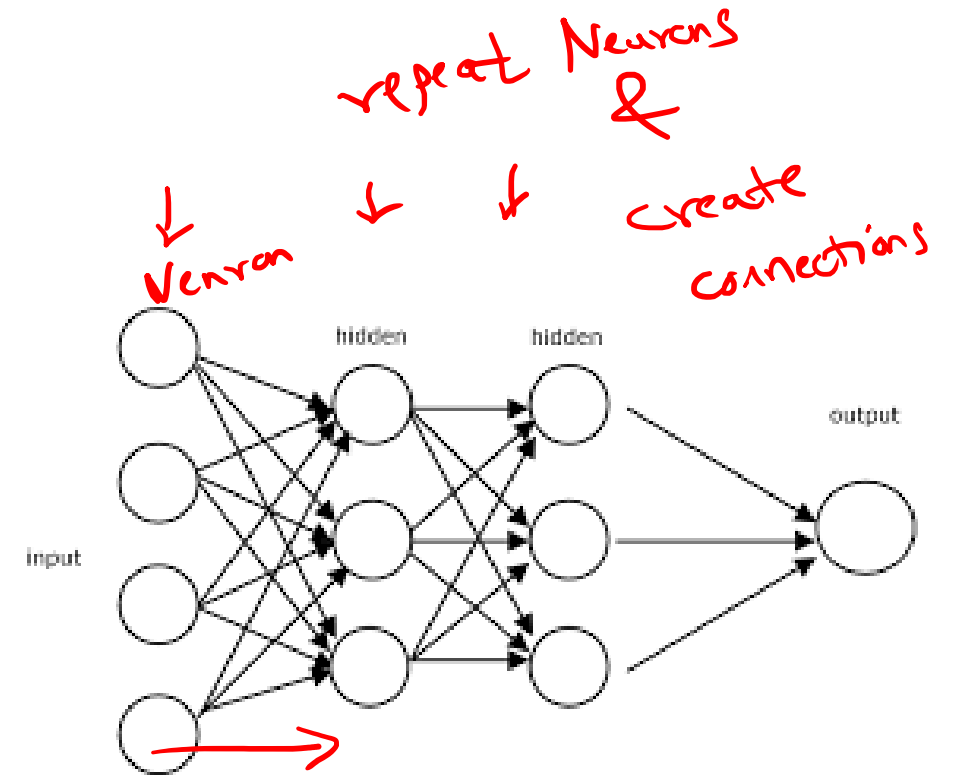
CNN (CV)  
~~RNN~~ (NLP)

Real brain: Wires can be criss-cross



# Forward Propagation

Model—used for making predictions

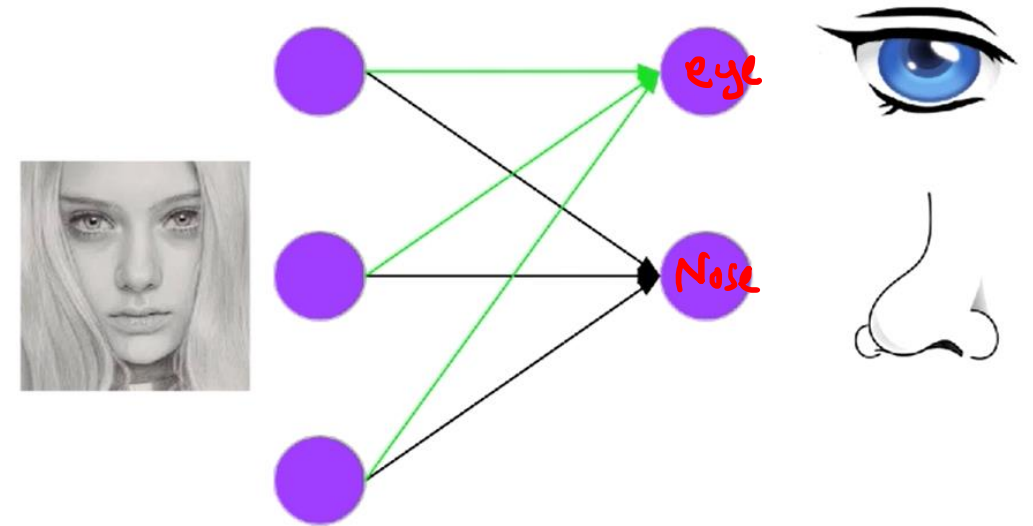


# Repeating the single neuron

- Each of these neurons may be calculating something different, via different weights

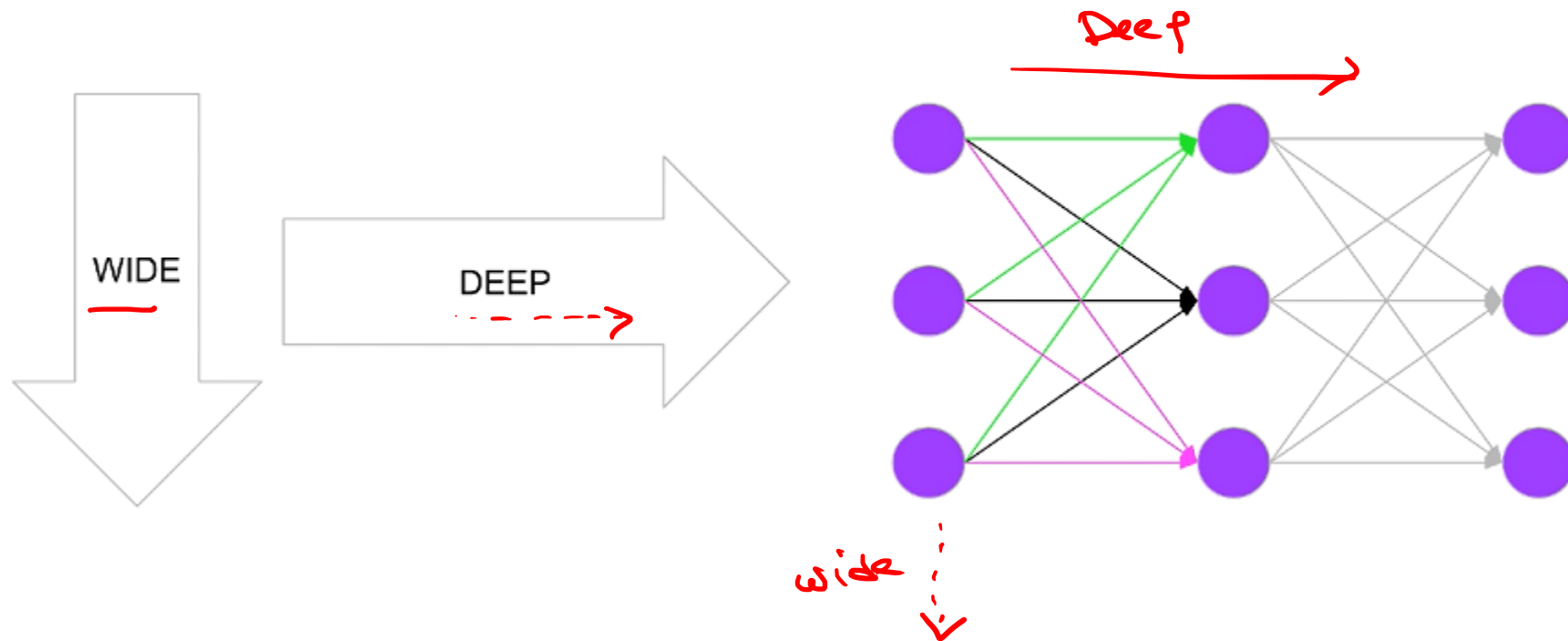
eg: If the i/p is Face

- (a) One neuron → EYE } different  
(b) One neuron → NOSE } features



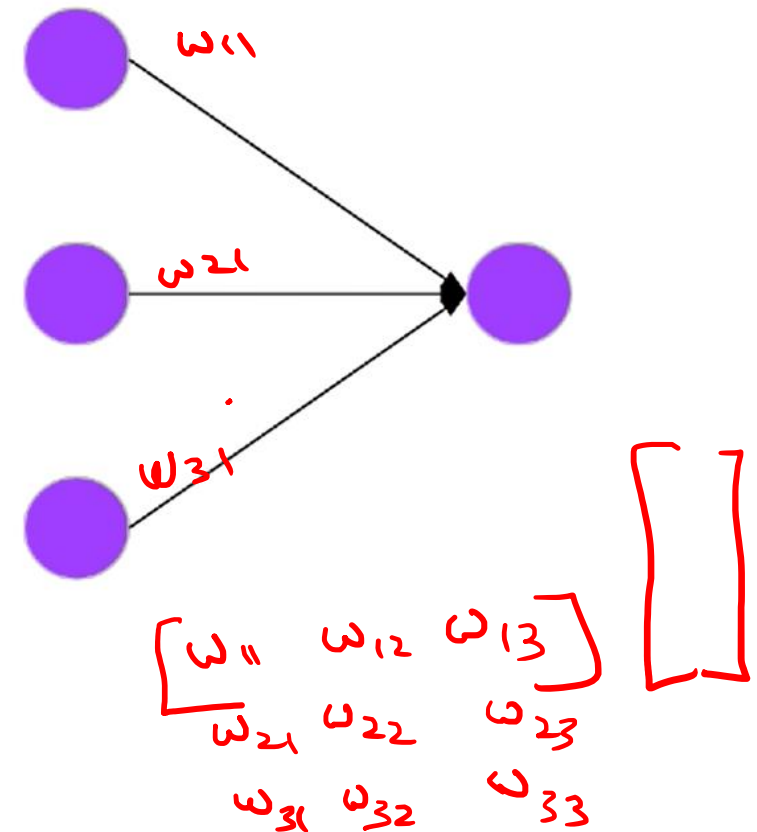
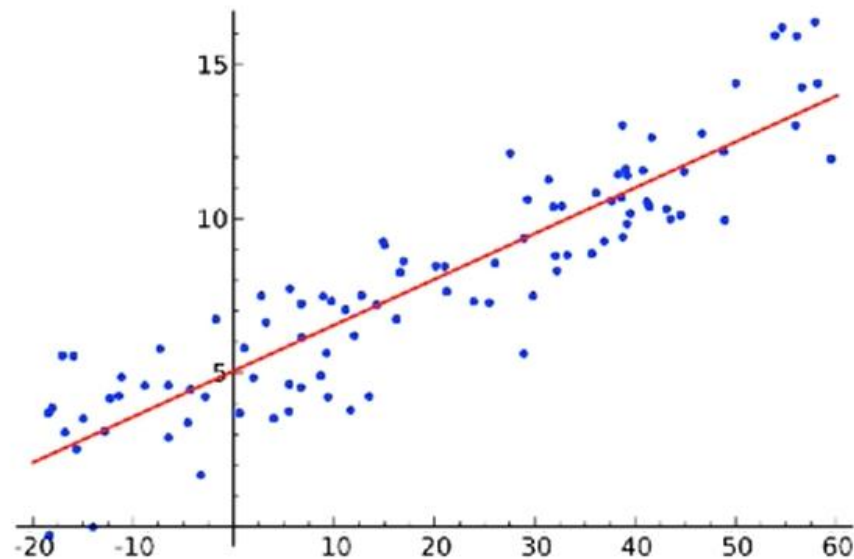
# Two important ways to extend a single neuron

1. The same inputs can be fed to multiple different neurons, each calculating something different (more neurons per layer) **WIDE**
2. Neurons in one layer can act as inputs to another layer **(DEEP)**





# Lines to Neurons



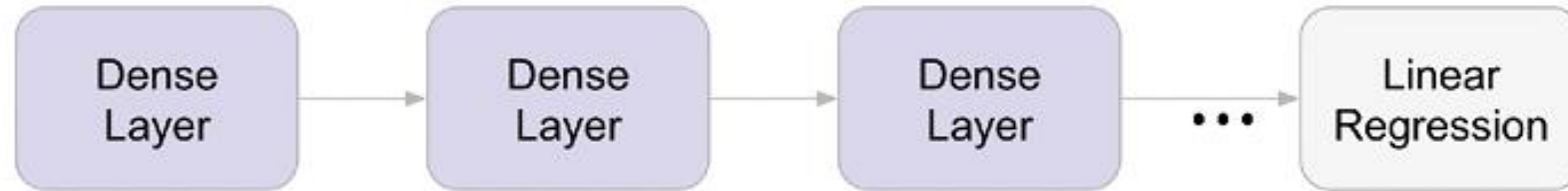
*A line :  $ax$  +  $b$*

*A neuron :  $\sigma(\underline{w}^T \underline{x} + \underline{b})$*

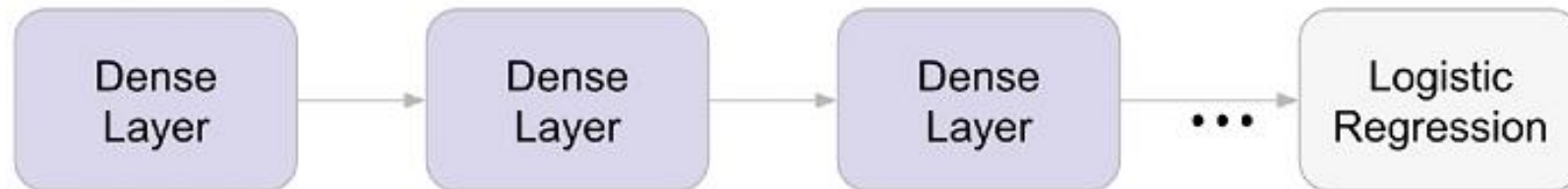
# Another perspective

- Each neural network is a feature transformation

Regression:



Classification:



# Feature hierarchies

- Researchers noticed that each layer learns increasingly complex features

Strokes - diff faces

0 0 0 0  
0 0 0 0  
8 8 8 8  
0 0 0 0

Eyes  
Nose  
Lips

in detail  
when we cover CNN

CNN  
how? 3

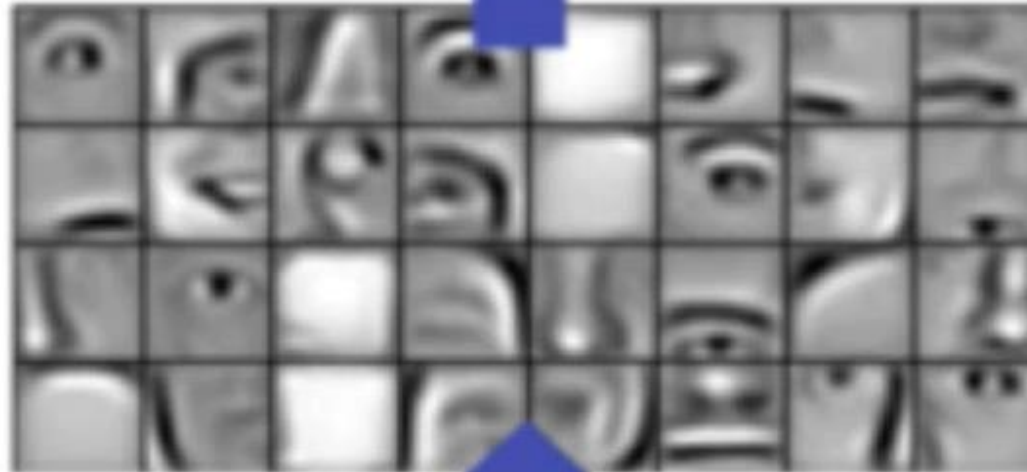
2

1



Different Face

Layer 3



Layer 2

Eye, Nose, Lips



Layer 1

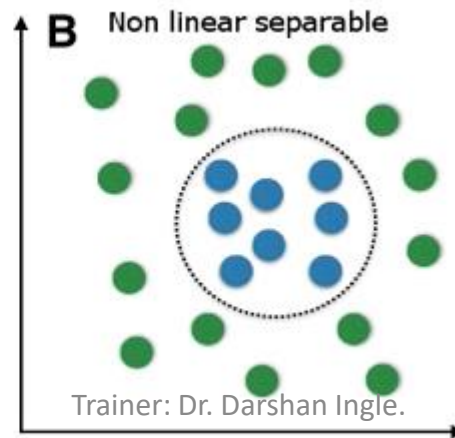
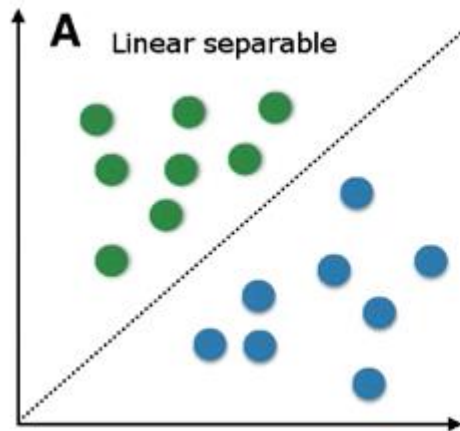
Strokes

# The Geometric Picture

- ML is nothing but a geometry problem
- Why are Neural Network so important?
- Why cant we just use a single neuron?
- The neuron is nice and interpretable
- Large weights = important feature
- Small weights = not important feature
- Unfortunately, the neuron (linear model) is not very expensive
- But true learning doesn't happen with a single neuron

# Making the line more complicated

- 2 ways to make our problem more complicated than “finding a line”
  1. Adding more input dimensions
  2. “Make the pattern non linear” (This is what we are concerned with now)



# TensorFlow Playground

- <https://playground.tensorflow.org/#activation=tanh&batchSize=10&dataset=circle&regDataset=reg-plane&learningRate=0.03&regularizationRate=0&noise=0&networkShape=4,2&seed=0.33964&showTestData=false&discretize=false&percTrainData=50&x=true&y=true&xTimesY=false&xSquared=false&ySquared=false&cosX=false&sinX=false&cosY=false&sinY=false&collectStats=false&problem=classification&initZero=false&hideText=false>

# Revisiting Activation Functions

Sigmoid: 0 and 1 (between)

$$\sigma(a) = \frac{1}{1+\exp(-a)}$$



mimics  
Biological Neuron



Problems?  $\therefore$  It is no longer used any more  
as exhaustively as it was used before.

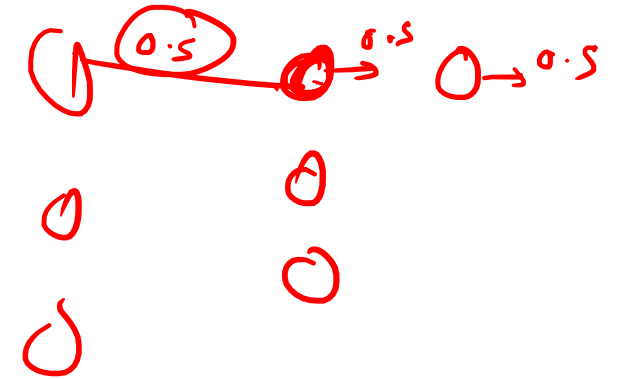
# Standardization

inp range (a) 1 to 5 million  
(b) 0 to 0.0001 } inp are not centered around 0

o/p of Sigmoid is b/w 0 & 1 & its middle is 0.5

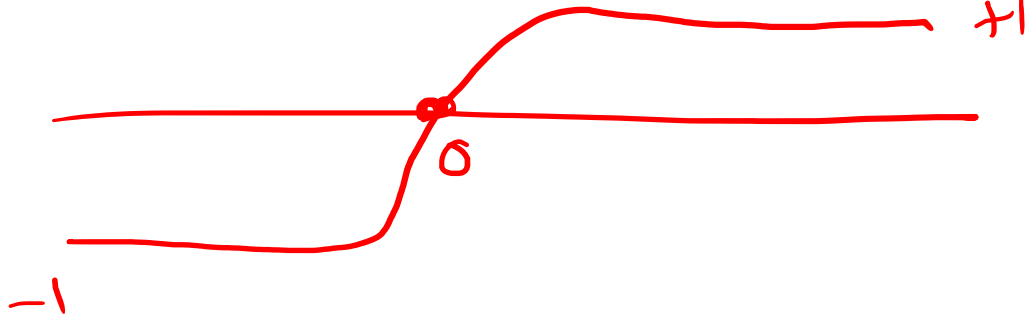
∴ o/p of Sigmoid can never be centered around 0.

Idea of Uniformity.





# Hyperbolic tangent (tanh)



$$\tanh(a) = \frac{\exp(2a)-1}{\exp(2a)+1} \quad \frac{\sinh}{\cosh}$$

# Still more problems

tanh is a little better than sigmoid, but still both are problematic.

None of researchers wanted to break away from the classical way of doing things.

Because of that, Prob! Vanishing Gradient Problem.

# Vanishing gradient problem

In 1980-1995, researchers were not able to create a Deep N.N.  
bcz they were using Sigmoid in each & every Neuron.

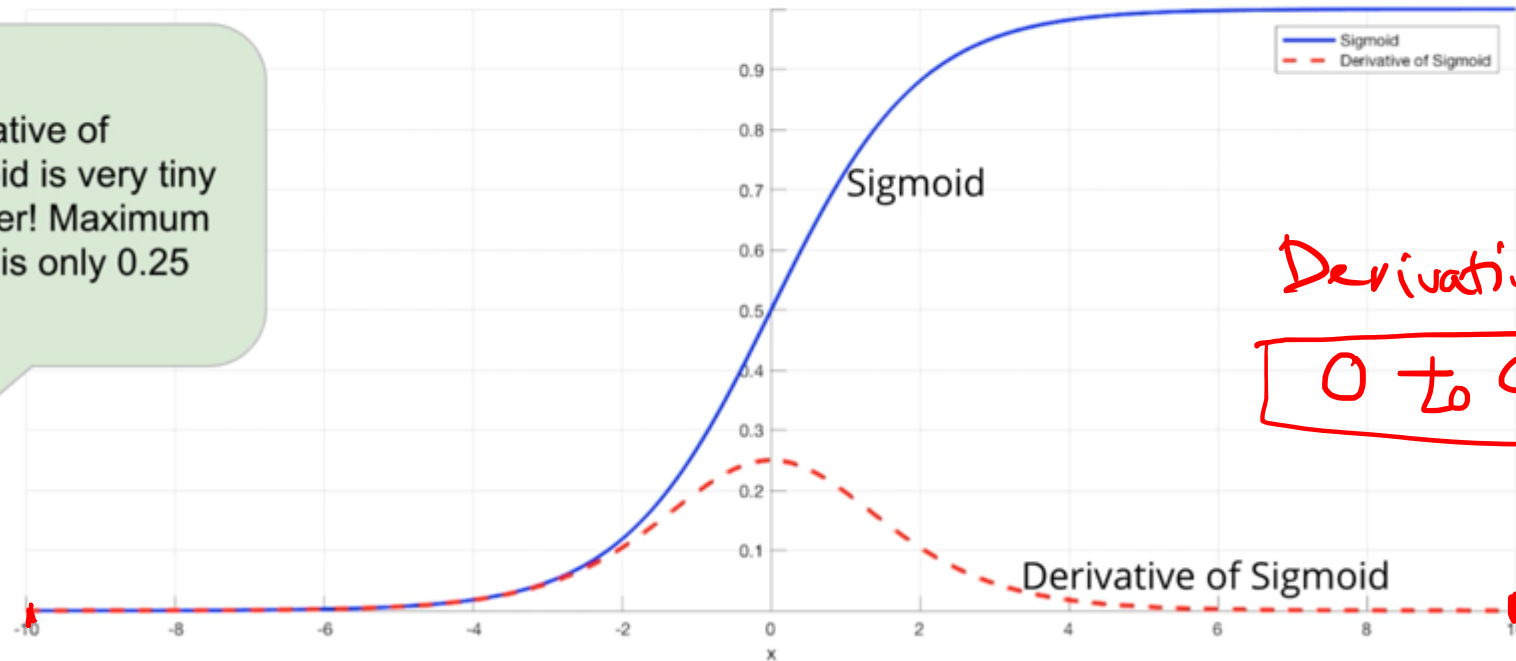
ReLU was not invented that time.

Bcz, everyone suffered from V.G. Problem,

# Vanishing gradient problem

Sigmoid  
$$w_{\text{new}} = w_{\text{old}} - \eta \cdot \frac{\partial L}{\partial w_{\text{old}}}$$
  
← Backprop. →

Derivative of sigmoid is very tiny number! Maximum value is only 0.25

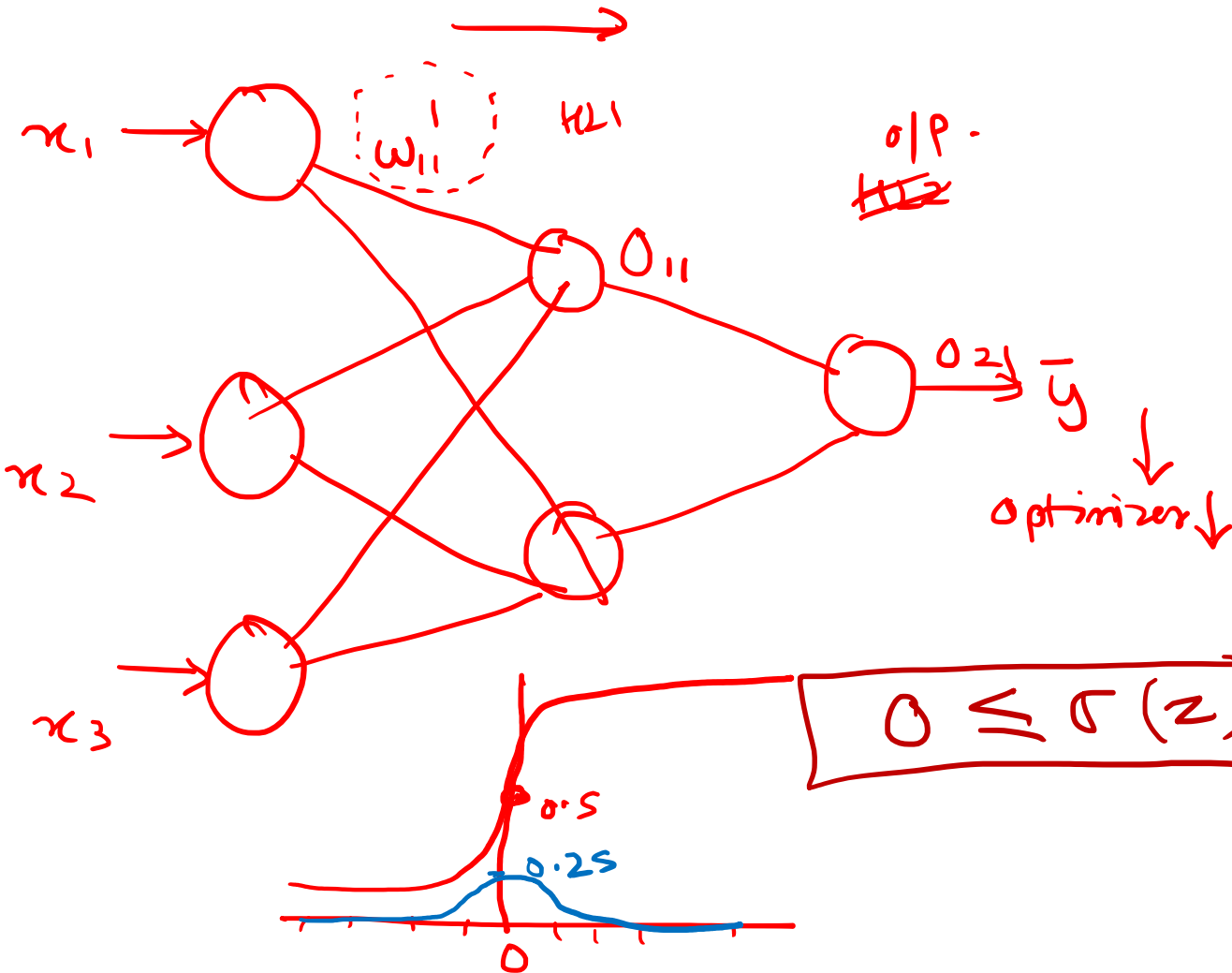


Derivative of sigmoid

0 to 0.25

# Vanishing gradient problem

$$\eta = 1$$



$$w'_{11_{\text{new}}} = w'_{11_{\text{old}}} - \eta \cdot \frac{\partial L}{\partial w'_{11_{\text{old}}}}$$

$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial o_{21}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w'_{11}}$$

(0 to 0.25) (0 to 0.25) = 0.5625

$$0 \leq \sigma(z) \leq 0.25$$

0.20	0.02	= 0.004
0.02	0.04	= 0.0008

$w'_{11_{\text{new}}} \approx w'_{11_{\text{old}}}$   
weight updation is not happening