



Inspire...Educate...Transform.

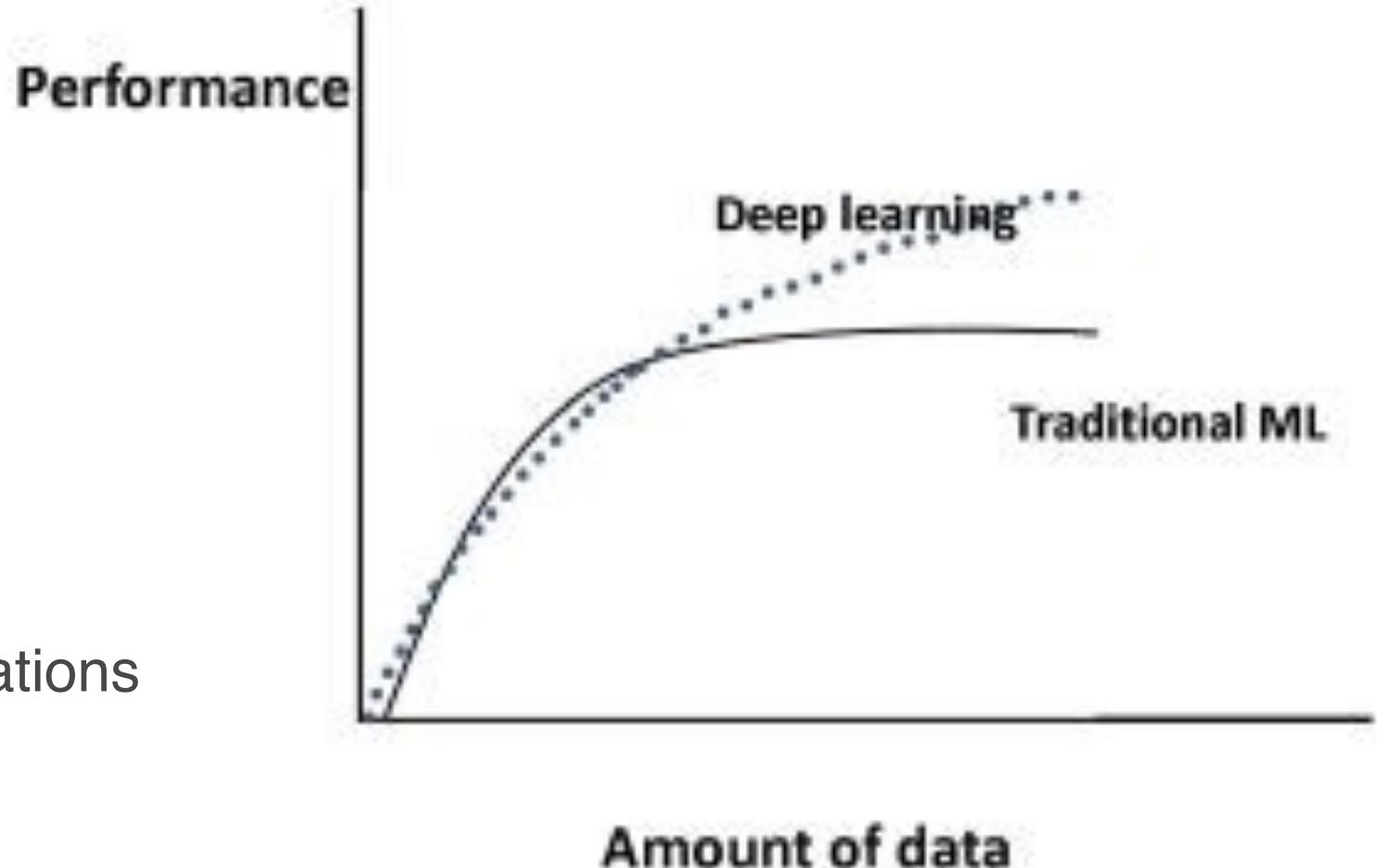
Artificial Neural Networks

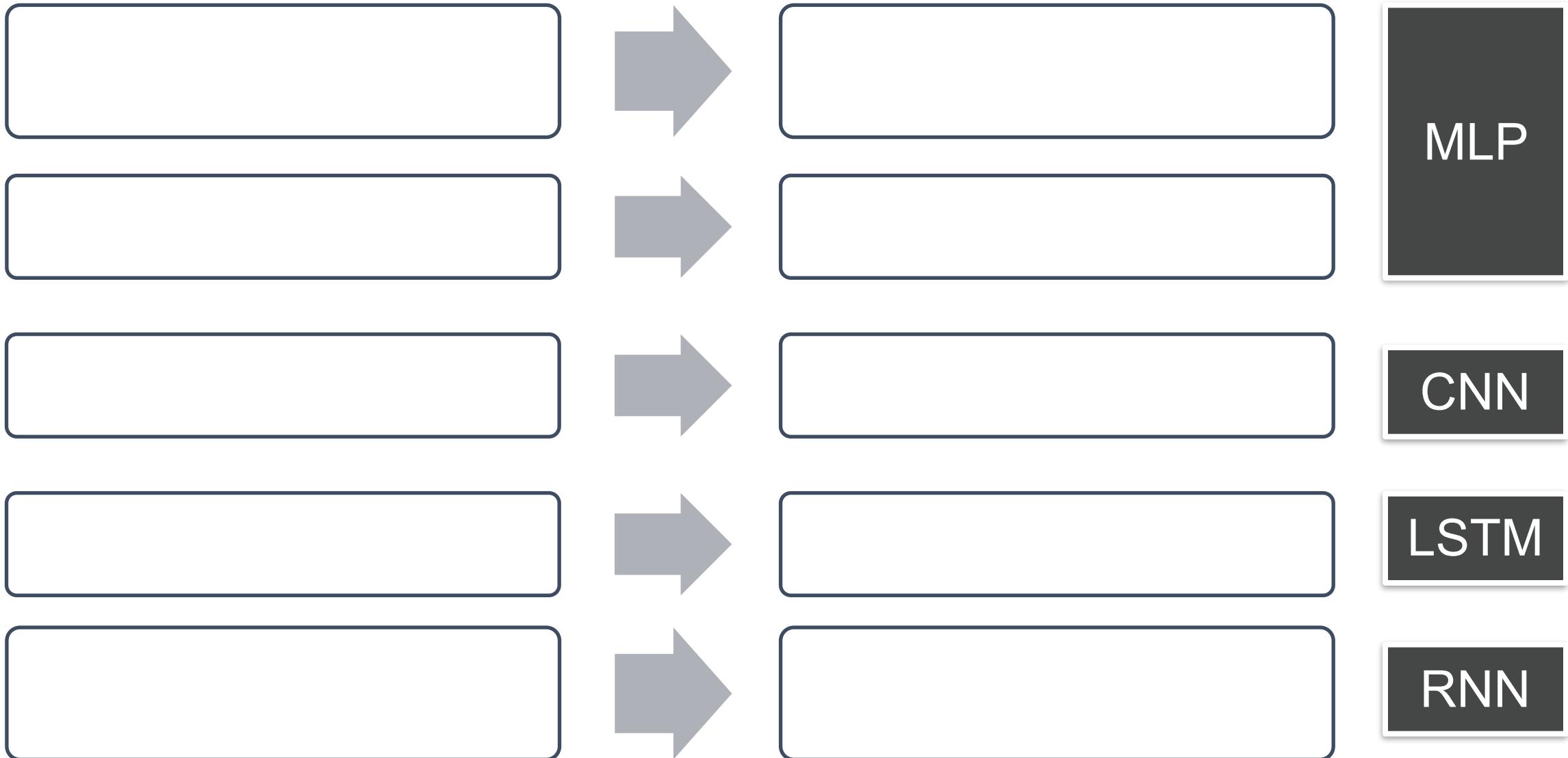
Parag Mantri, PhD

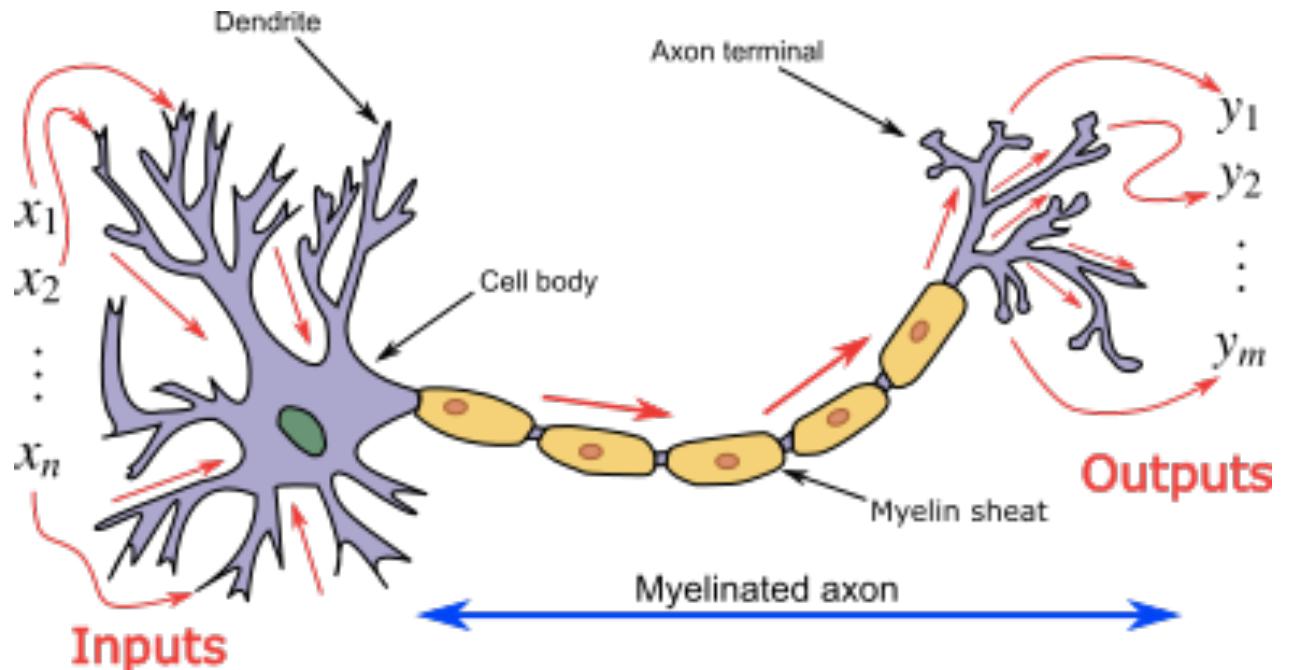
Neural Networks

- Basic idea behind ANN and why we need it now
- Understanding of single neuron or perceptron
- XOR example of why linear classification does not work
- Intuition behind need for activation functions and types. (RELU, Sigmoid, etc)
- Gradient descent review
- Equations of loss and cost function.
- Basic idea of forward and backward propagation
- Mathematical details behind backward propagation

- Data
 - Cell phones
 - Sensors
 - Camera
 - Wearables
 - Digitalization
- Computer Power
 - Smaller and faster
- Algorithms
 - Problem Specific innovations



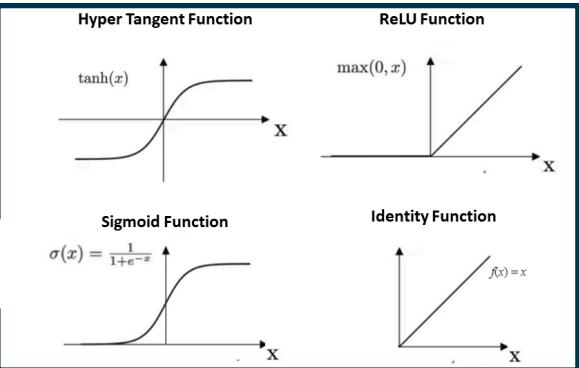
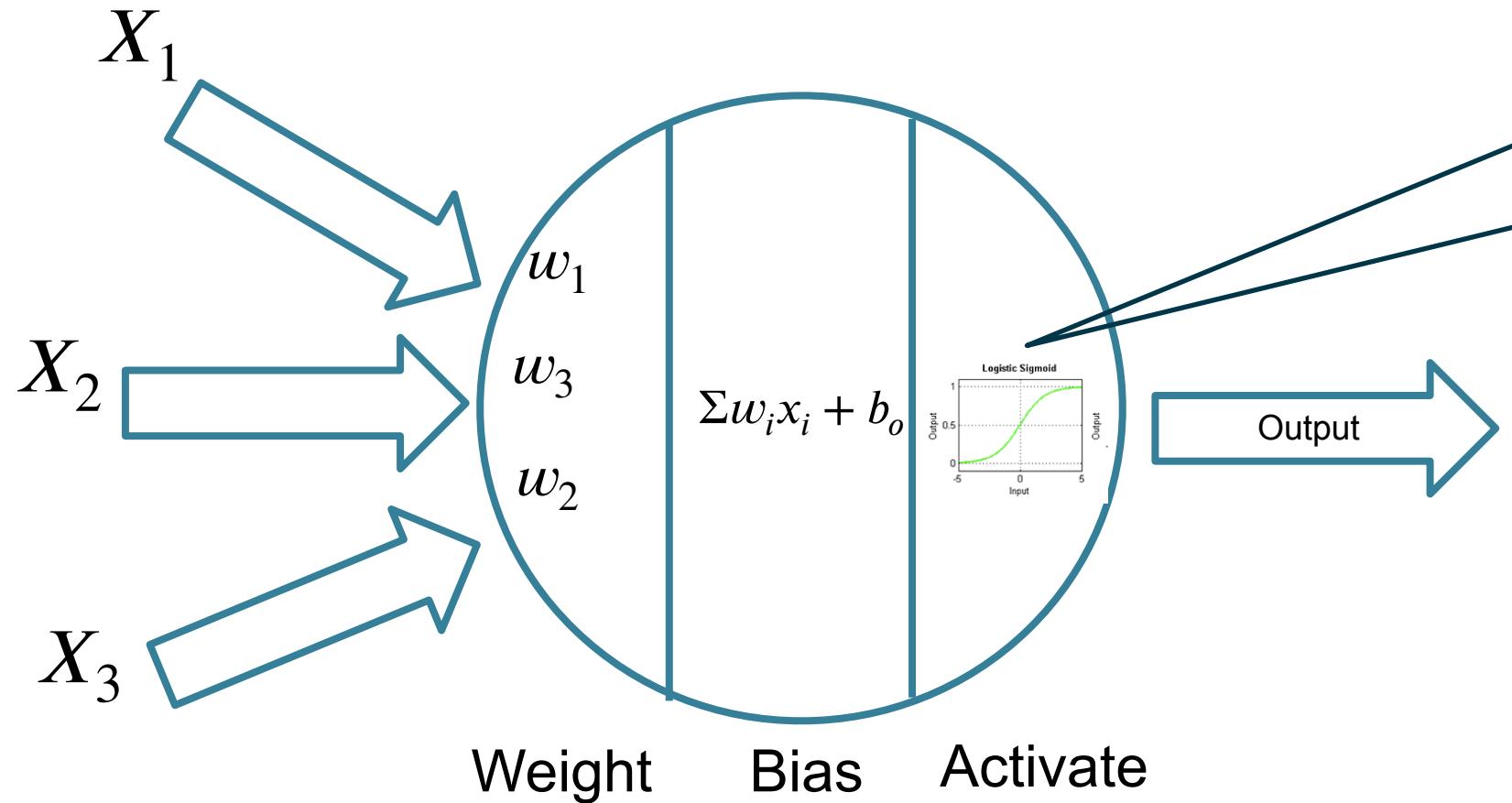




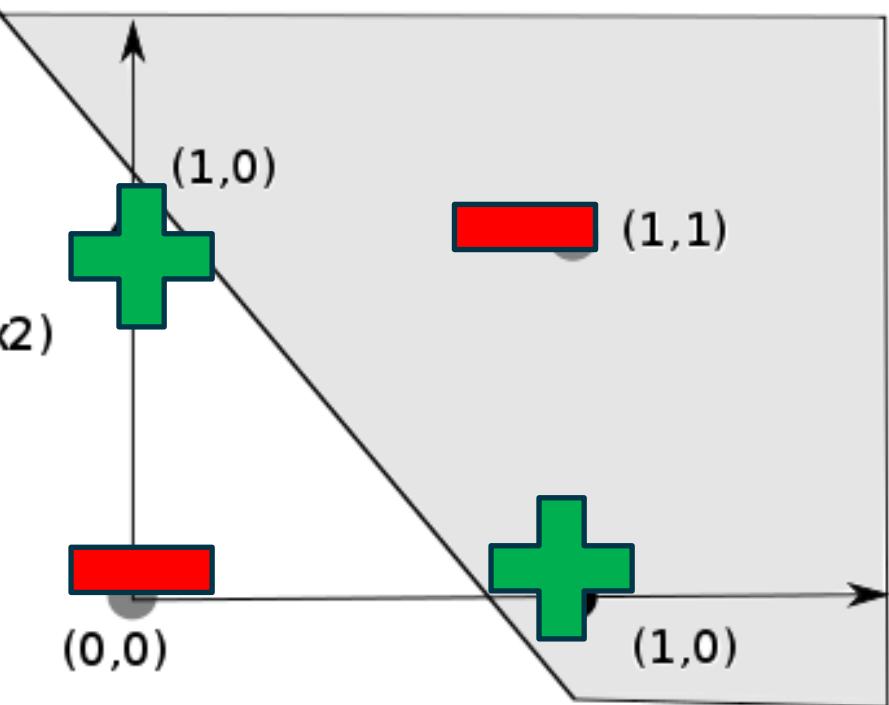
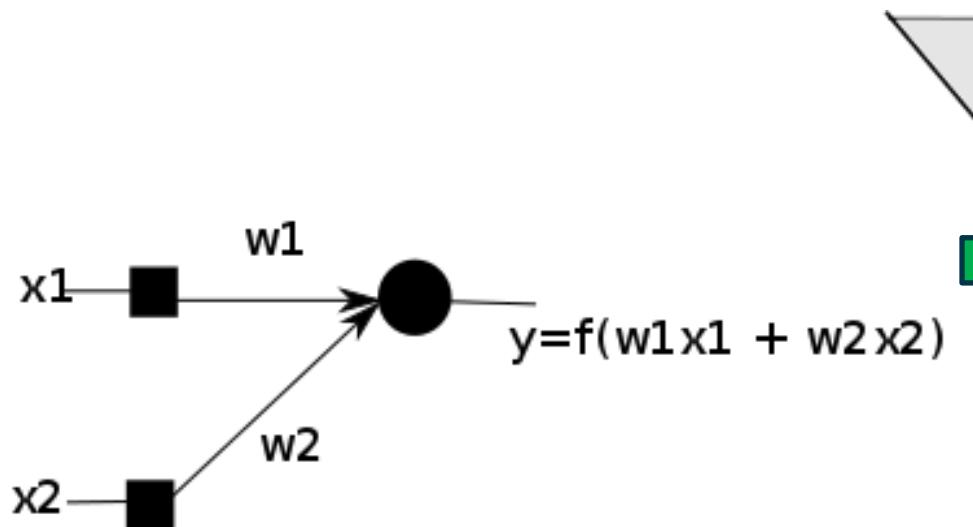
Goal Understand the brain

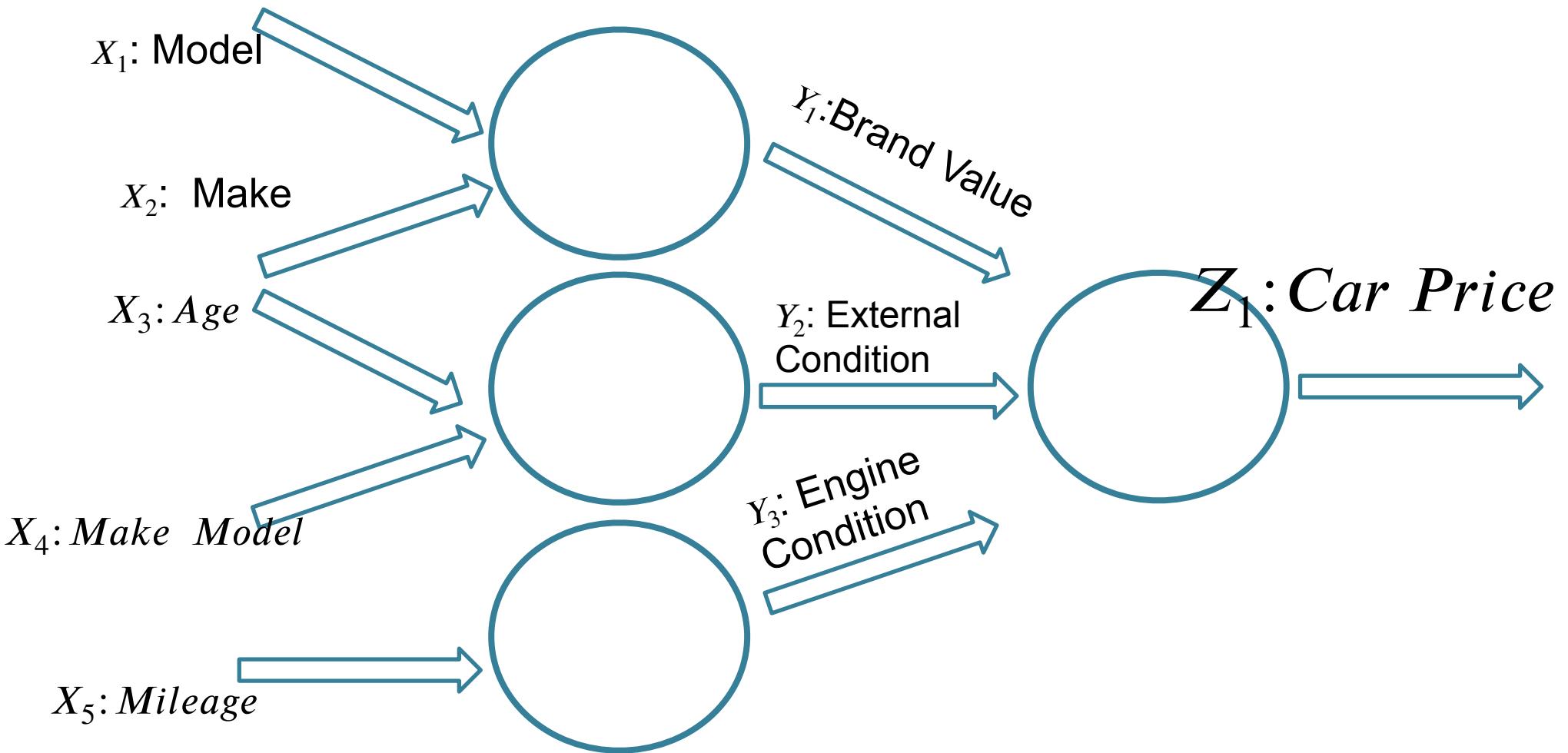
Emulate the Brain

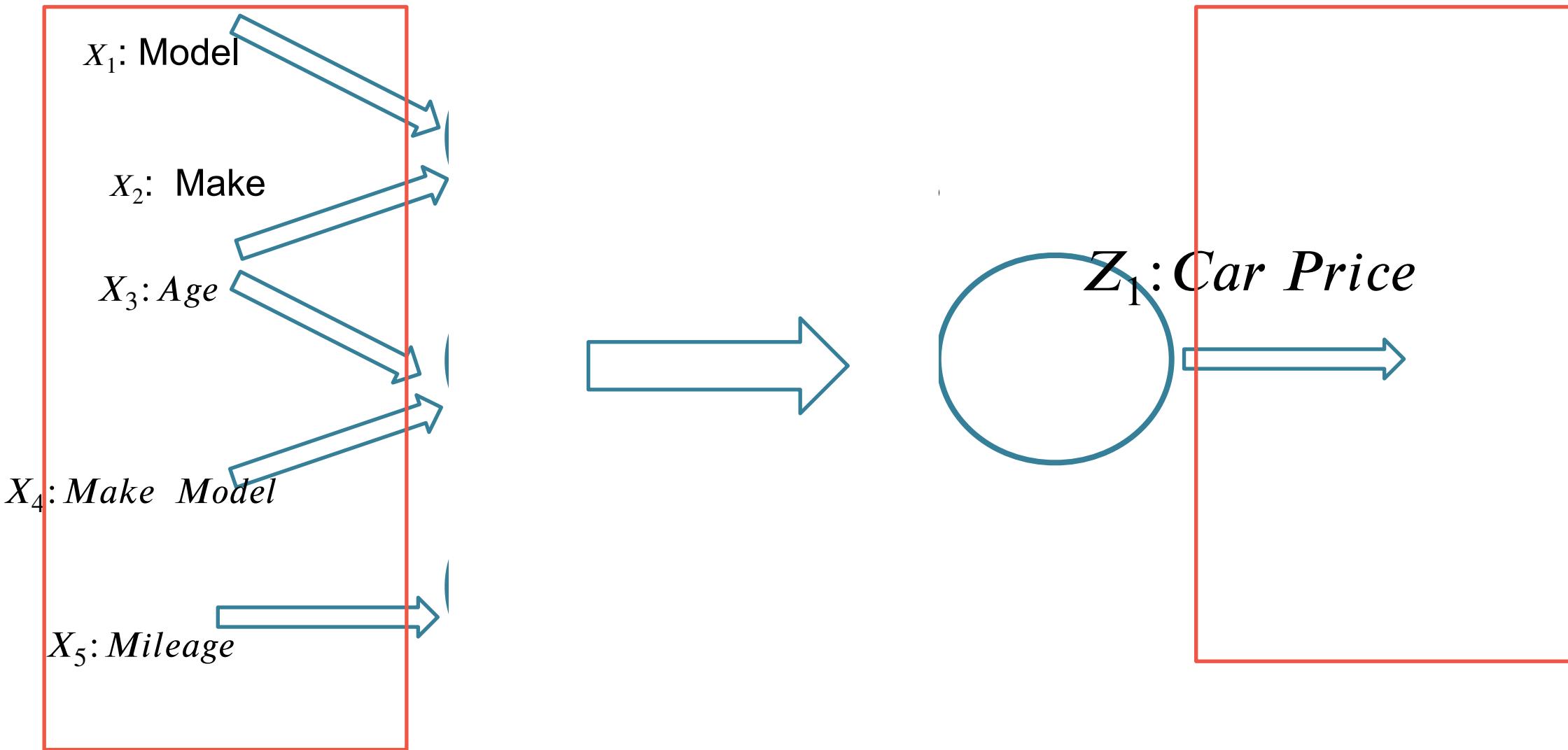
ANN Did we achieve the goal?



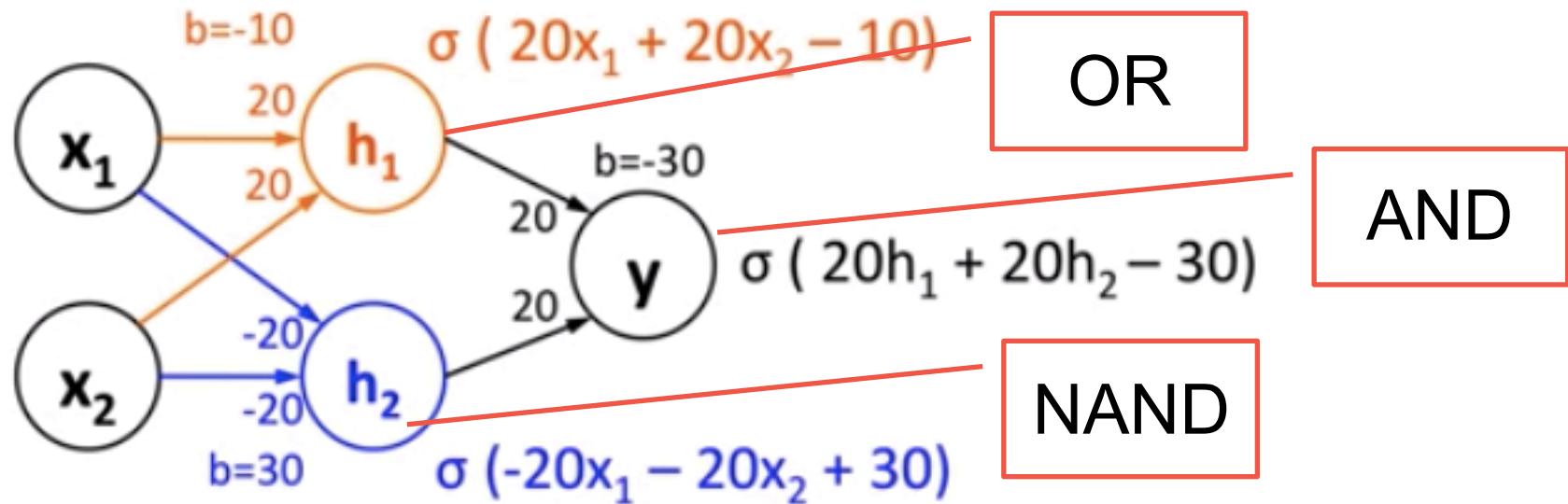
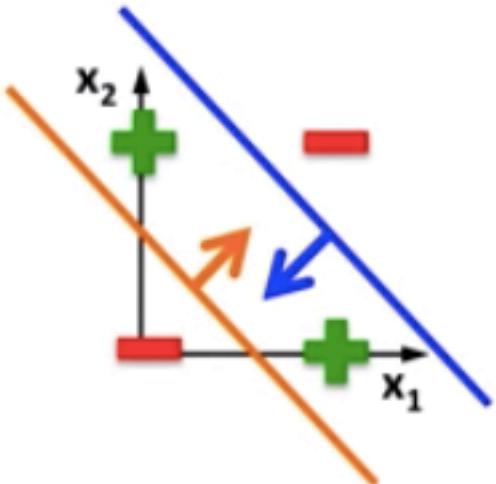
x1	x2	y
0	0	0
1	0	1
0	1	1
1	1	0







Linear classifiers
cannot solve this



$$\sigma(20*0 + 20*0 - 10) \approx 0$$

$$\sigma(20*1 + 20*1 - 10) \approx 1$$

$$\sigma(20*0 + 20*1 - 10) \approx 1$$

$$\sigma(20*1 + 20*0 - 10) \approx 1$$

$$\sigma(-20*0 - 20*0 + 30) \approx 1$$

$$\sigma(-20*1 - 20*1 + 30) \approx 0$$

$$\sigma(-20*0 - 20*1 + 30) \approx 1$$

$$\sigma(-20*1 - 20*0 + 30) \approx 1$$

$$\sigma(20*0 + 20*1 - 30) \approx 0$$

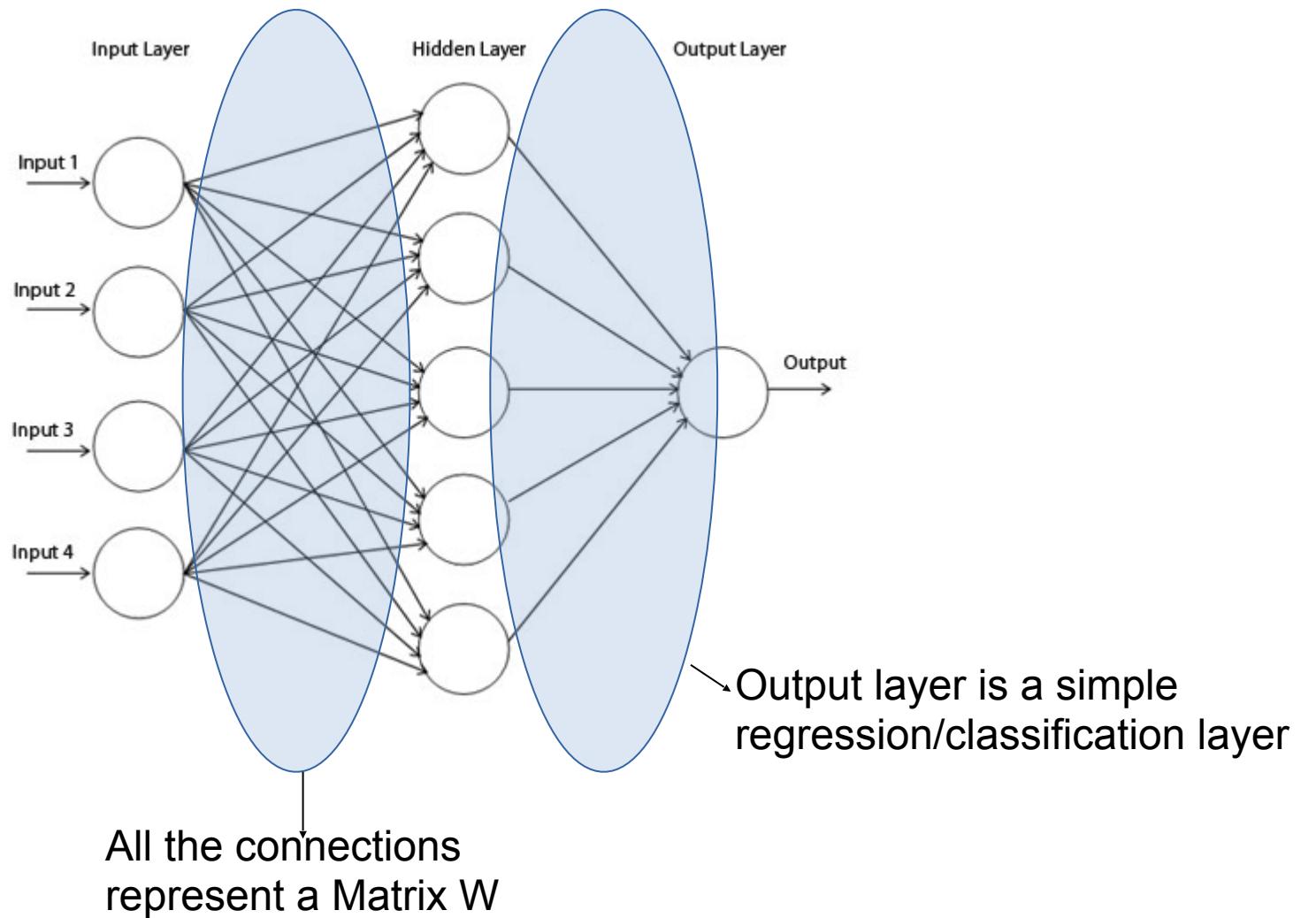
$$\sigma(20*1 + 20*0 - 30) \approx 0$$

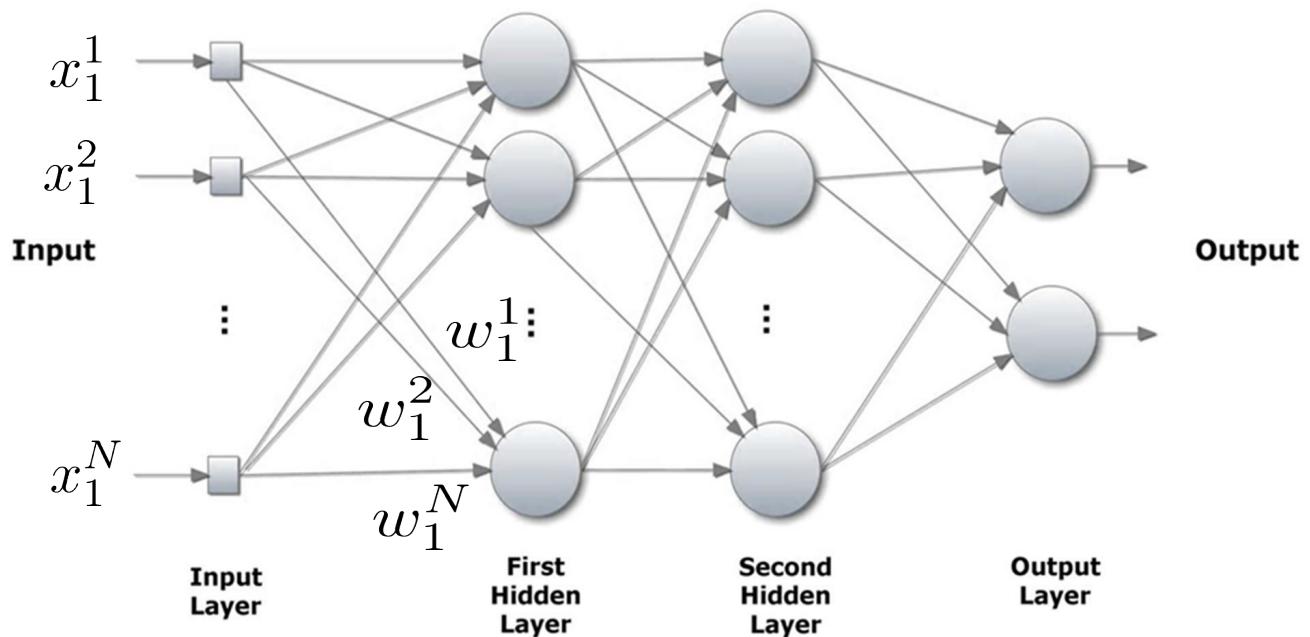
$$\sigma(20*1 + 20*1 - 30) \approx 1$$

$$\sigma(20*1 + 20*1 - 30) \approx 1$$

Each layer is a transformation into new space

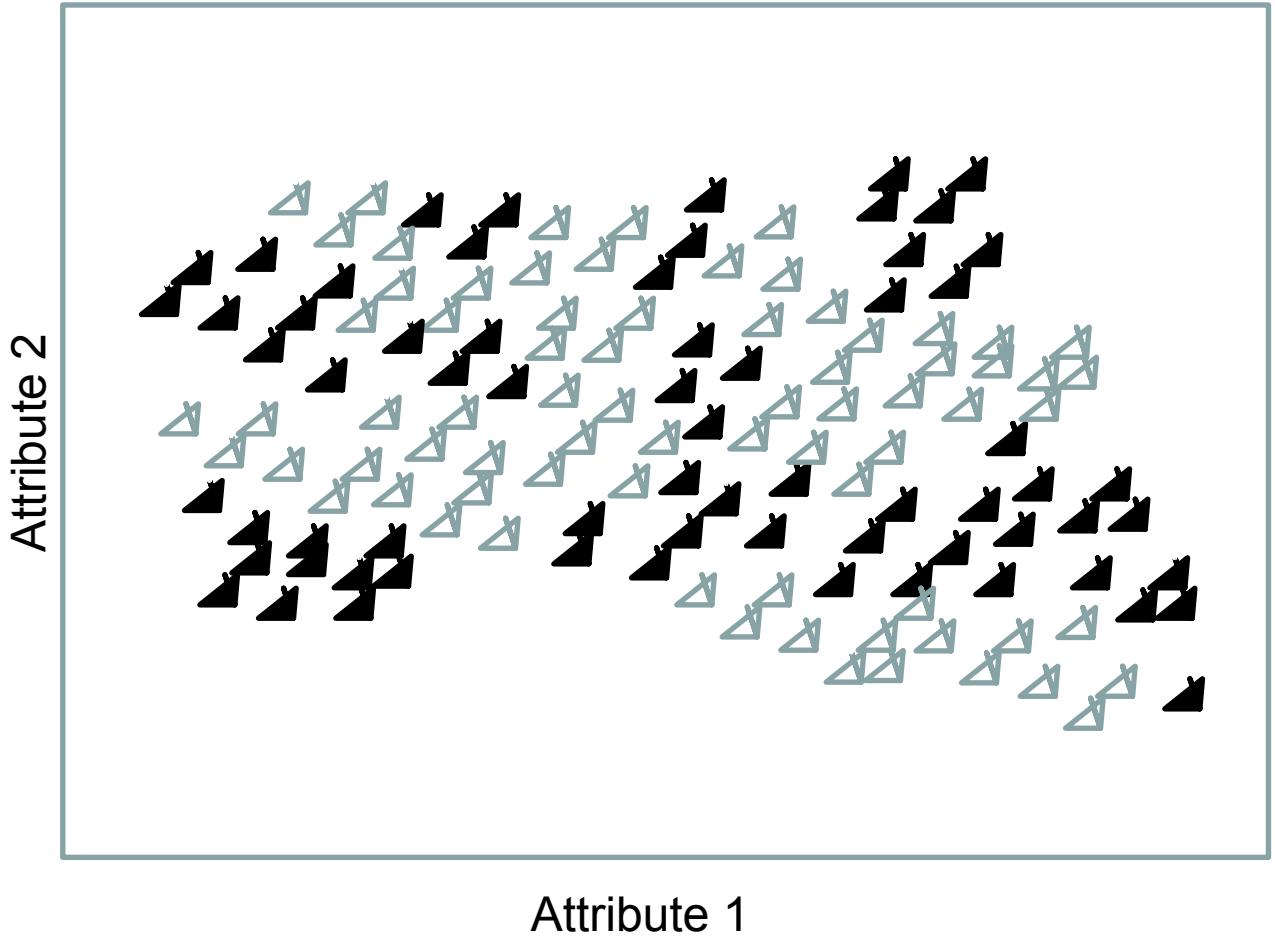
Ref: [Neural Networks Video](#) by Victor Lavrenko

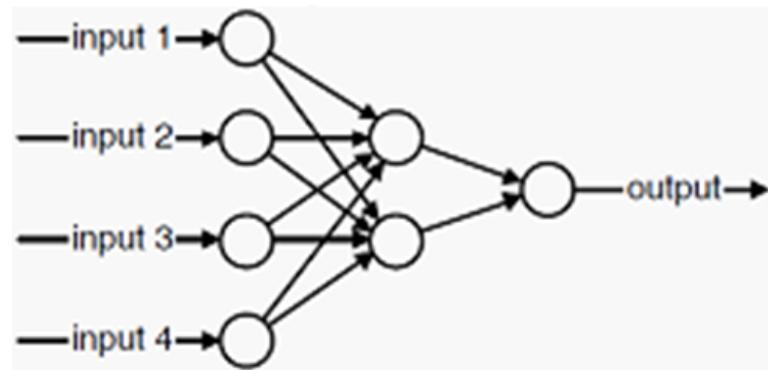
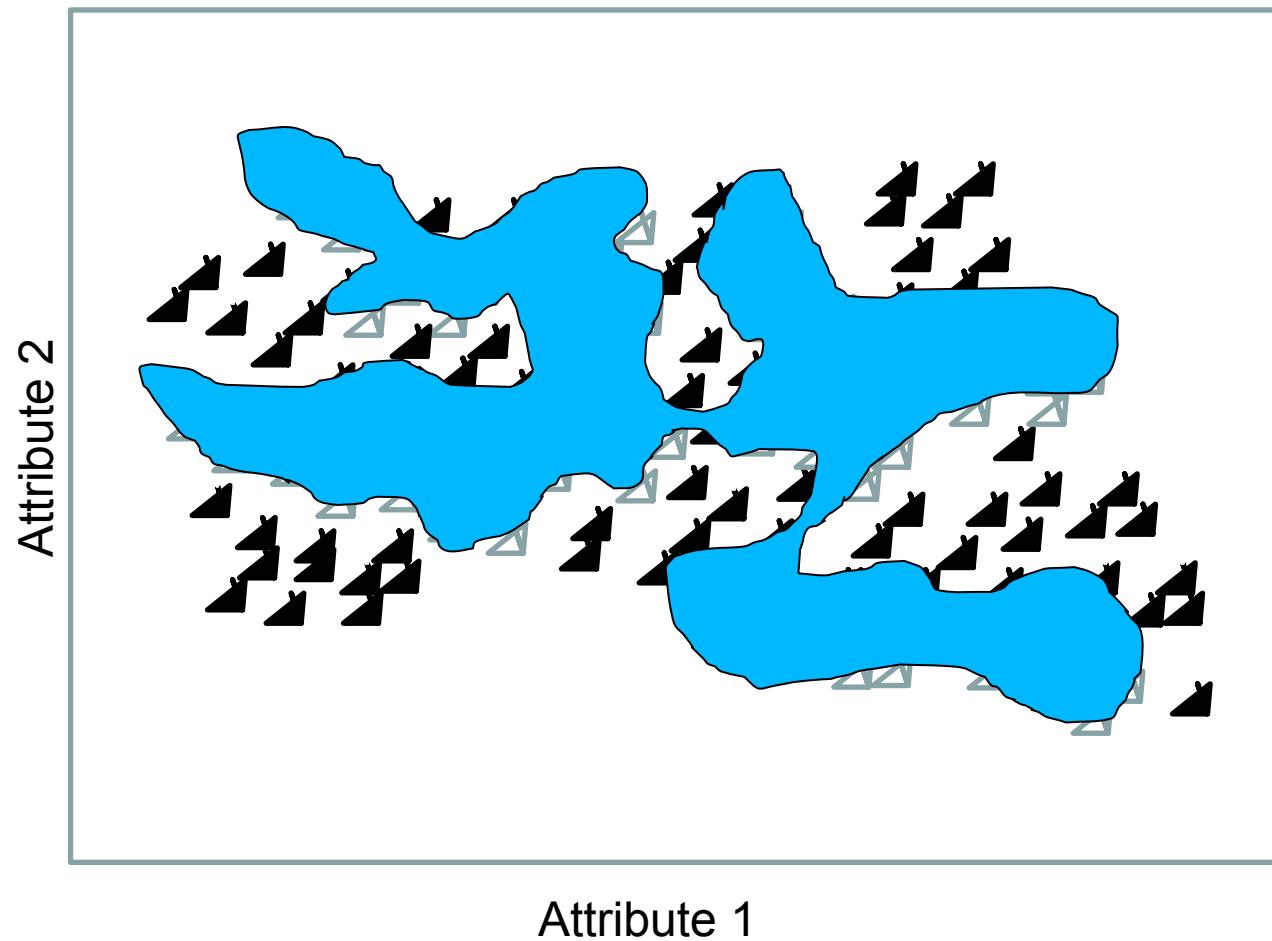




A multilayer perceptron

- Feed forward network
- Fully connected



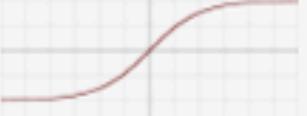
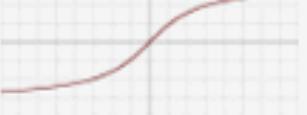


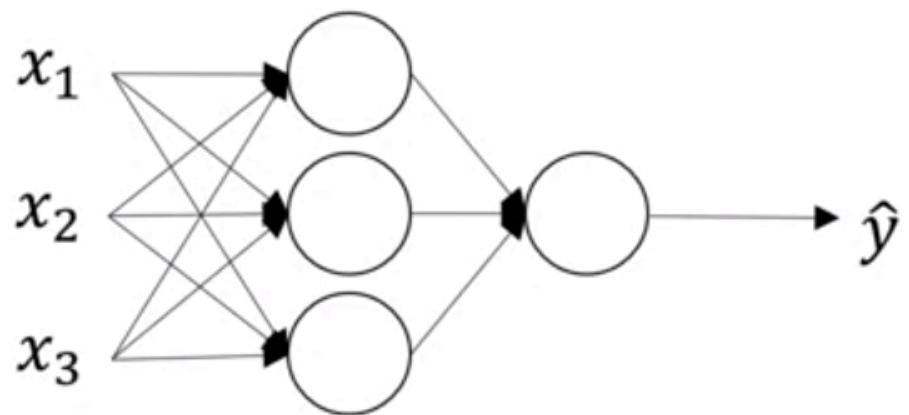
- 90%: One hidden layer
- 10%: Two hidden layers
- Rarely more than that! It might lead to overfitting!

- Most structured-data problems can be handled with at-most 2 layers. Rare are problems where higher layers can be justified.
- “rules of thumb” exist but are absolute rules. Experiment and validate and then cross validate like any other hyper parameter

See ftp://ftp.sas.com/pub/neural/FAQ3.html#A_hu

Activation Functions

Name	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) \stackrel{?}{=} \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a. k. a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
Rectified Linear Unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$



Given x :

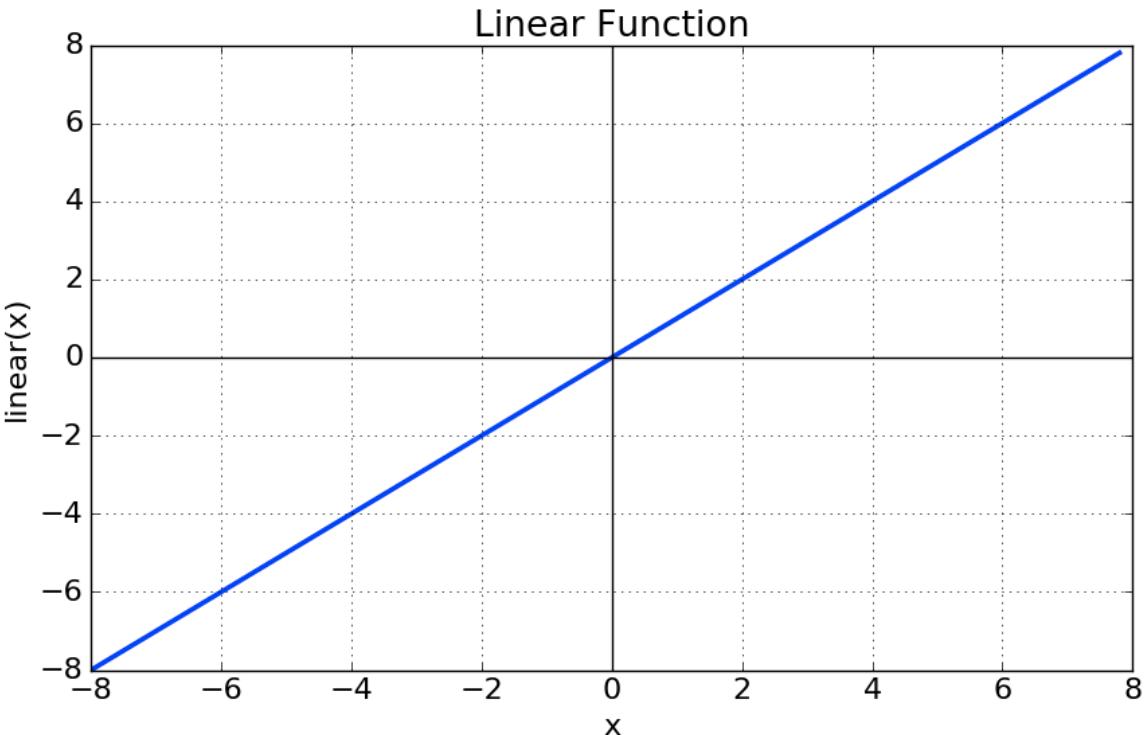
$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

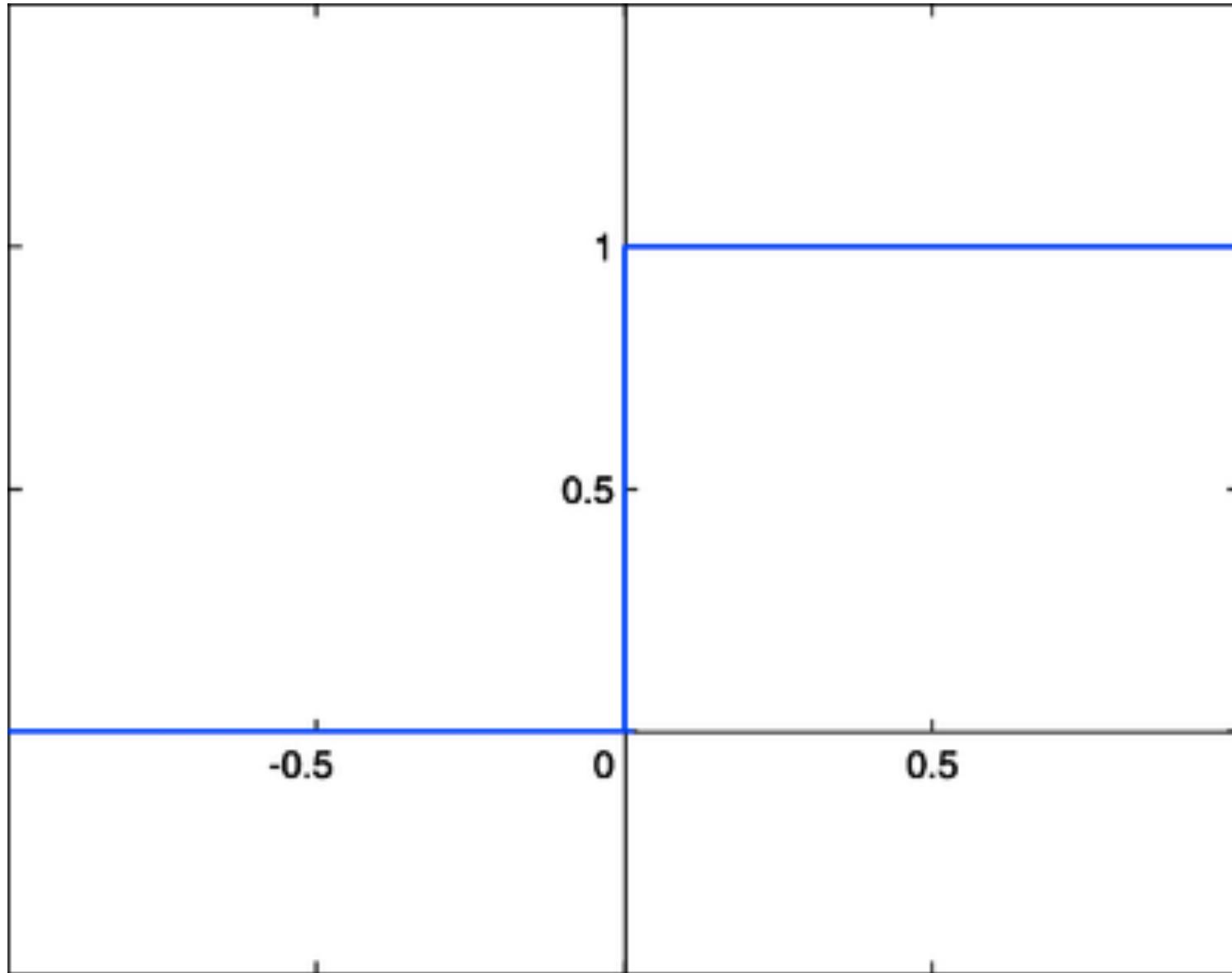
$$a^{[2]} = g^{[2]}(z^{[2]})$$

Activation Function

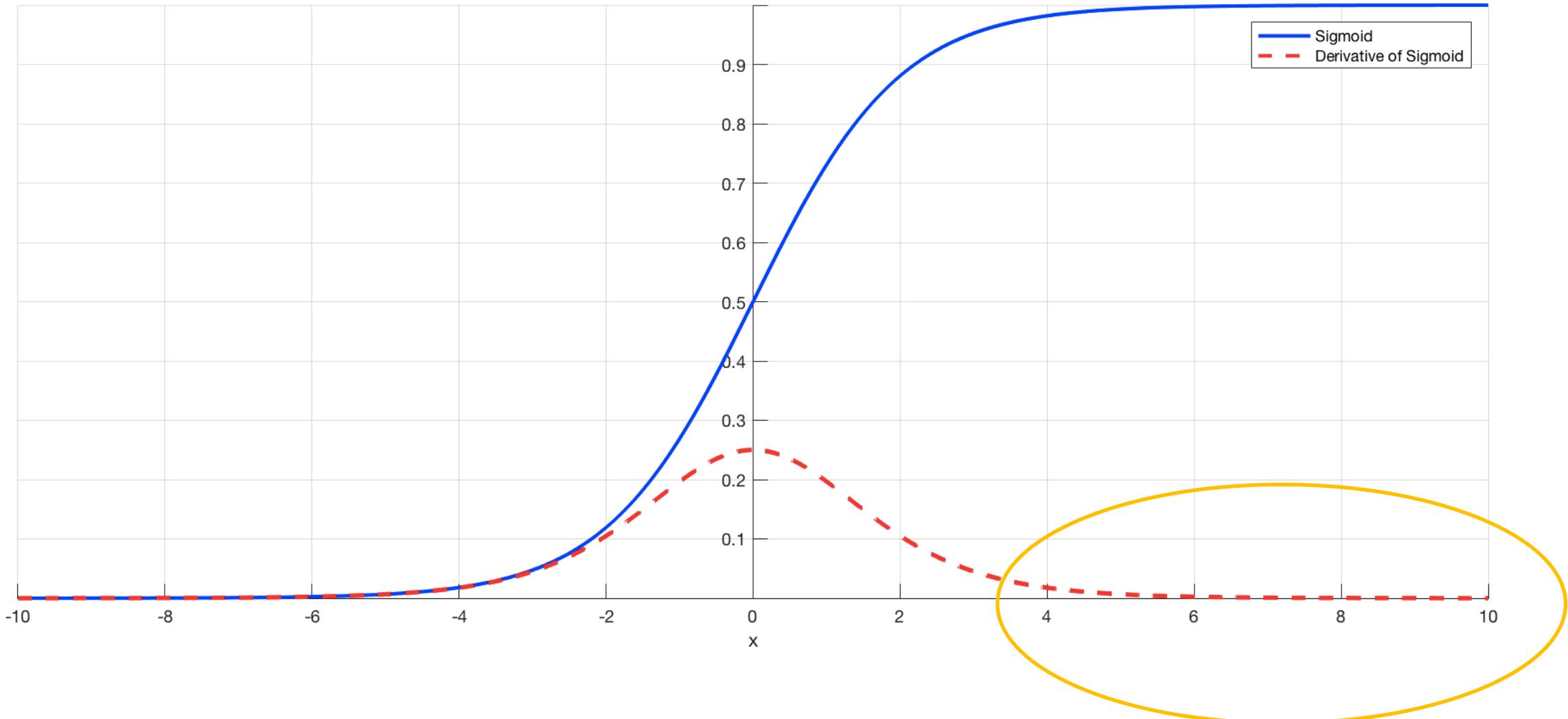


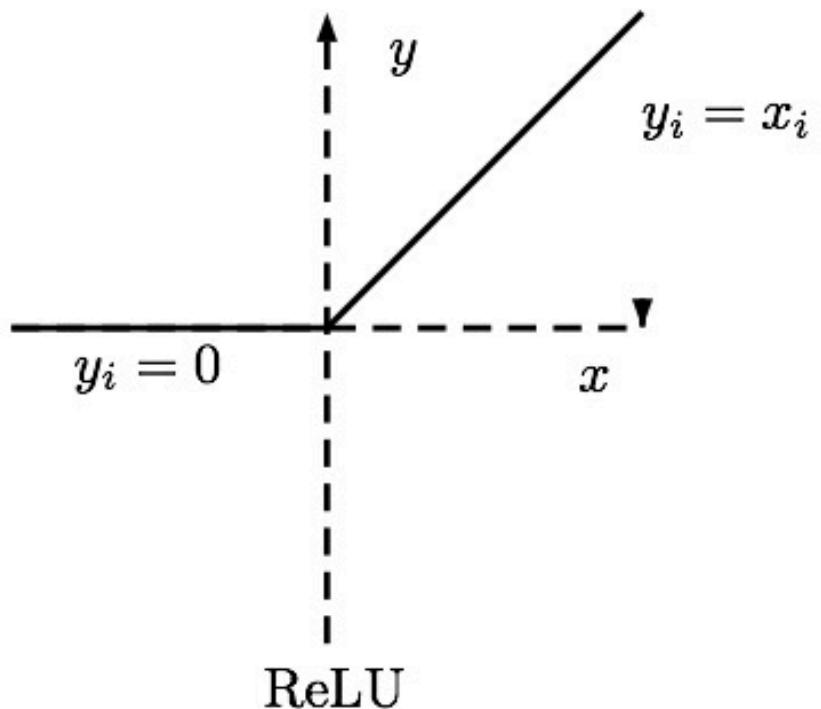
Q: What will be $Z[2]$ if we have linear activation function ?

- Linear Activation Function Output
 - $W^{[2]}W^{[1]}X + W^{[2]}b^{[1]} + b^{[2]}$
- Linear activation function defeats the purpose of multiple layer
- i.e. All the layers can be expressed as single neuron

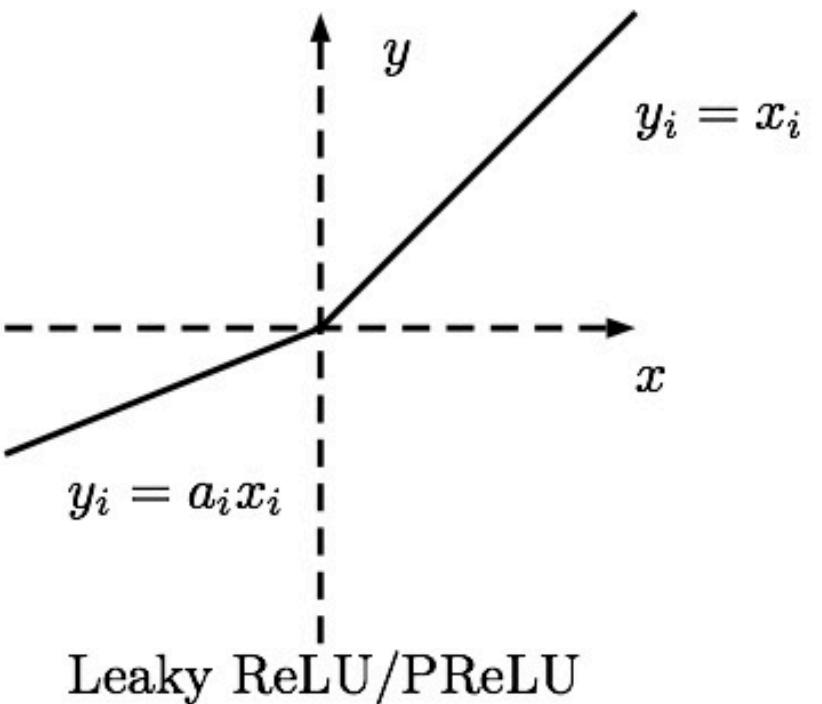


- Does not represent derivative learning.
- Suitable for output layer only



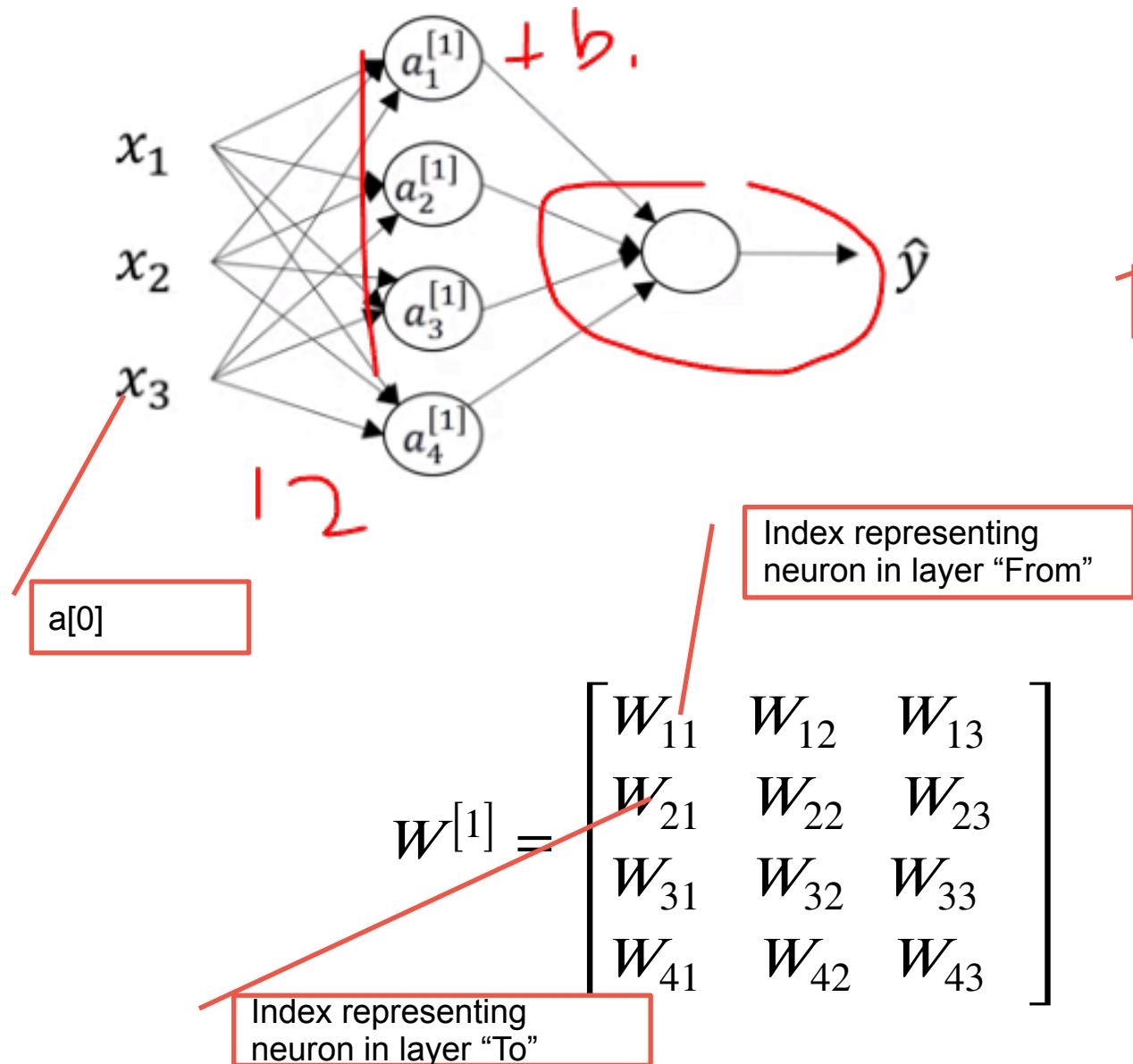


$$f(X) = \max((0, x))$$



Linear Algebra Convention

- n: Number of Features
- m: Number of training samples
- $X \in \mathbb{R}^n$, $y \in (0,1)$. (Classification Problem)
- Shape of X: n rows by m columns (What is the shape of a data frame?)
- Shape of y: 1 row by m columns



Given input x :

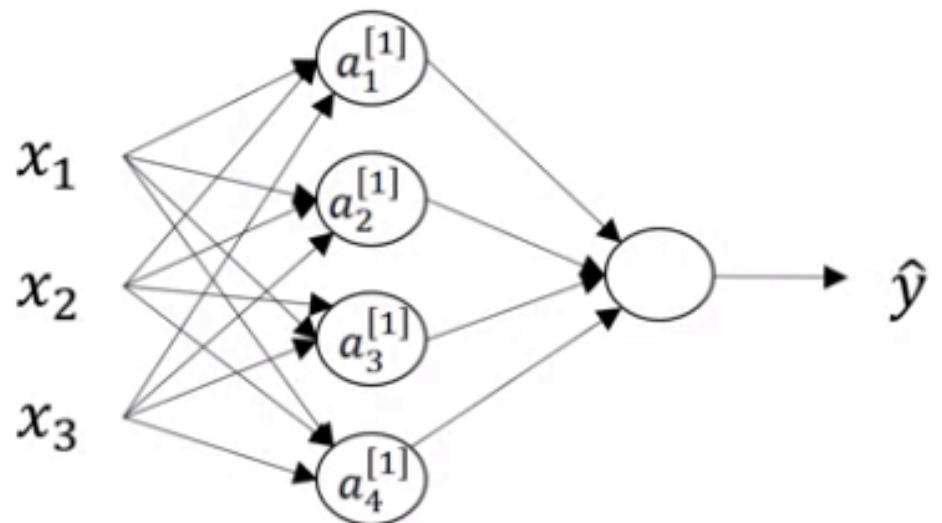
$$\begin{aligned} z^{[1]} &= W^{[1]}x + b^{[1]} && \text{Size } 4 \times 1 \\ a^{[1]} &= \sigma(z^{[1]}) && \text{Size } 4 \times 3 \end{aligned}$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

$$z = \sigma(x + b)$$

Layer Number: $Z^{(2)}$
 Neuron Number: Z_1



Given input x :

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

What will be the shape of all the elements in this layer ?

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2.$$

Desired
Output

NN Output

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)],$$

hidden

hidden

logits

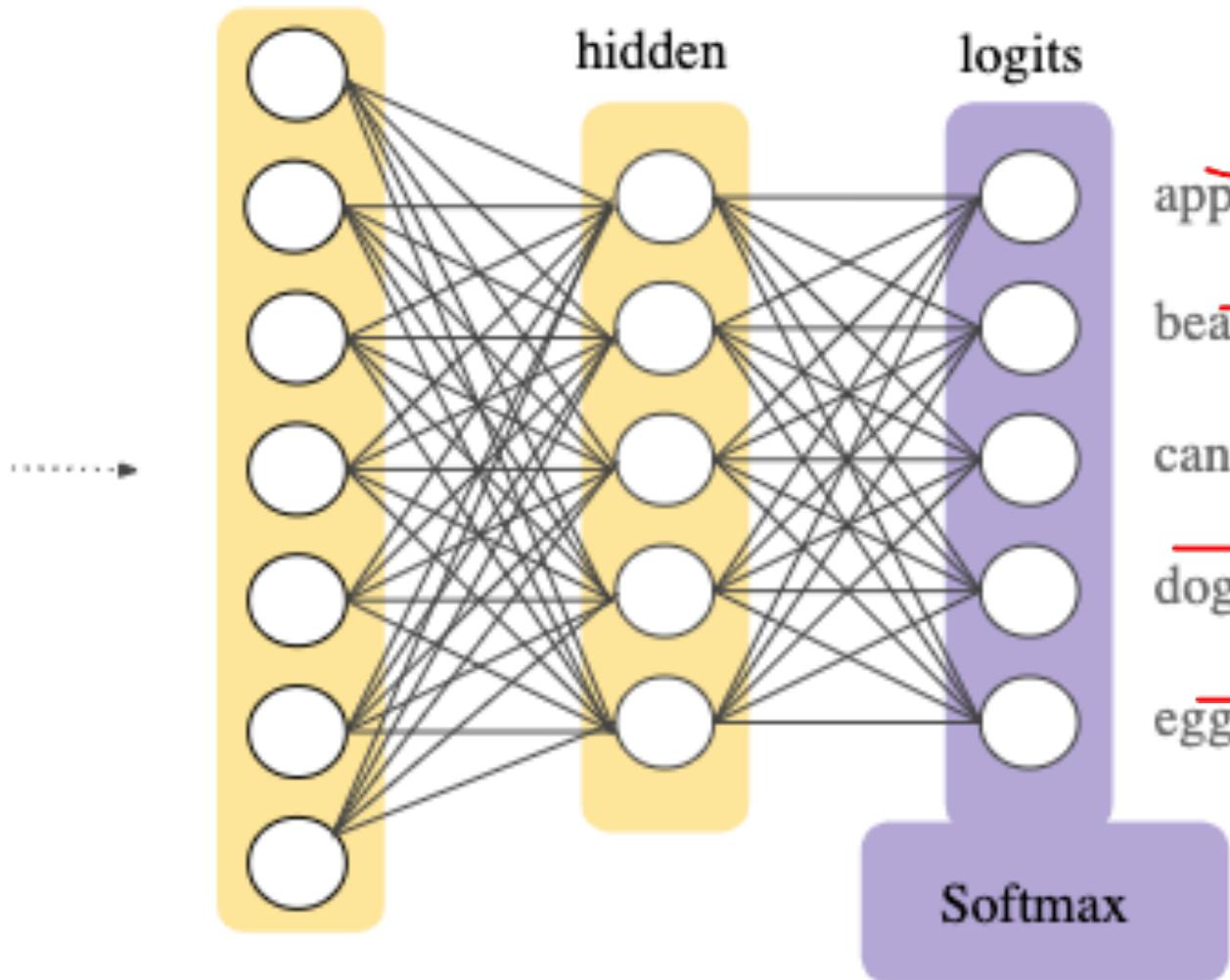
apple: yes/no?

bear: yes/no?

candy: yes/no?

dog: yes/no?

egg: yes/no?



Gradient Descent

$\hat{y} = \sigma(w^T x + b)$ where $\sigma(z) = \frac{1}{1 + e^{-z}}$

- Loss function (Individual training example)

$$L = y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

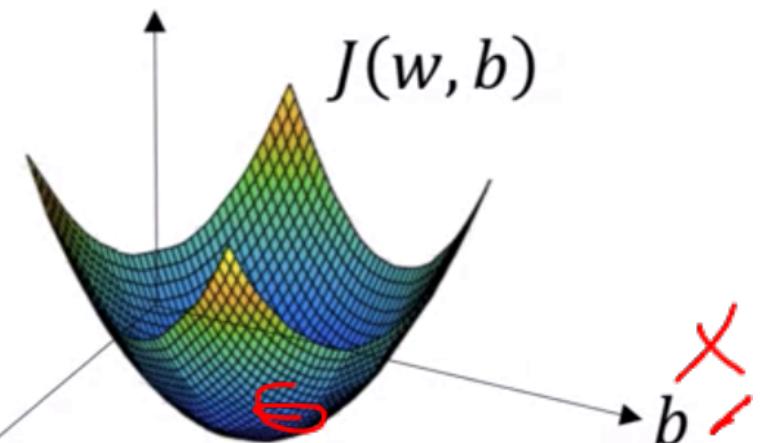
- Cost function (over entire data set)

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

$$\hat{y} = \begin{cases} 1 & \rightarrow 1 \\ 0 & \rightarrow 0 \end{cases}$$

$$L = \begin{cases} 0 & \hat{y} = 1 \\ 1 & \hat{y} = 0 \end{cases}$$

W is actually an
dimensional vector



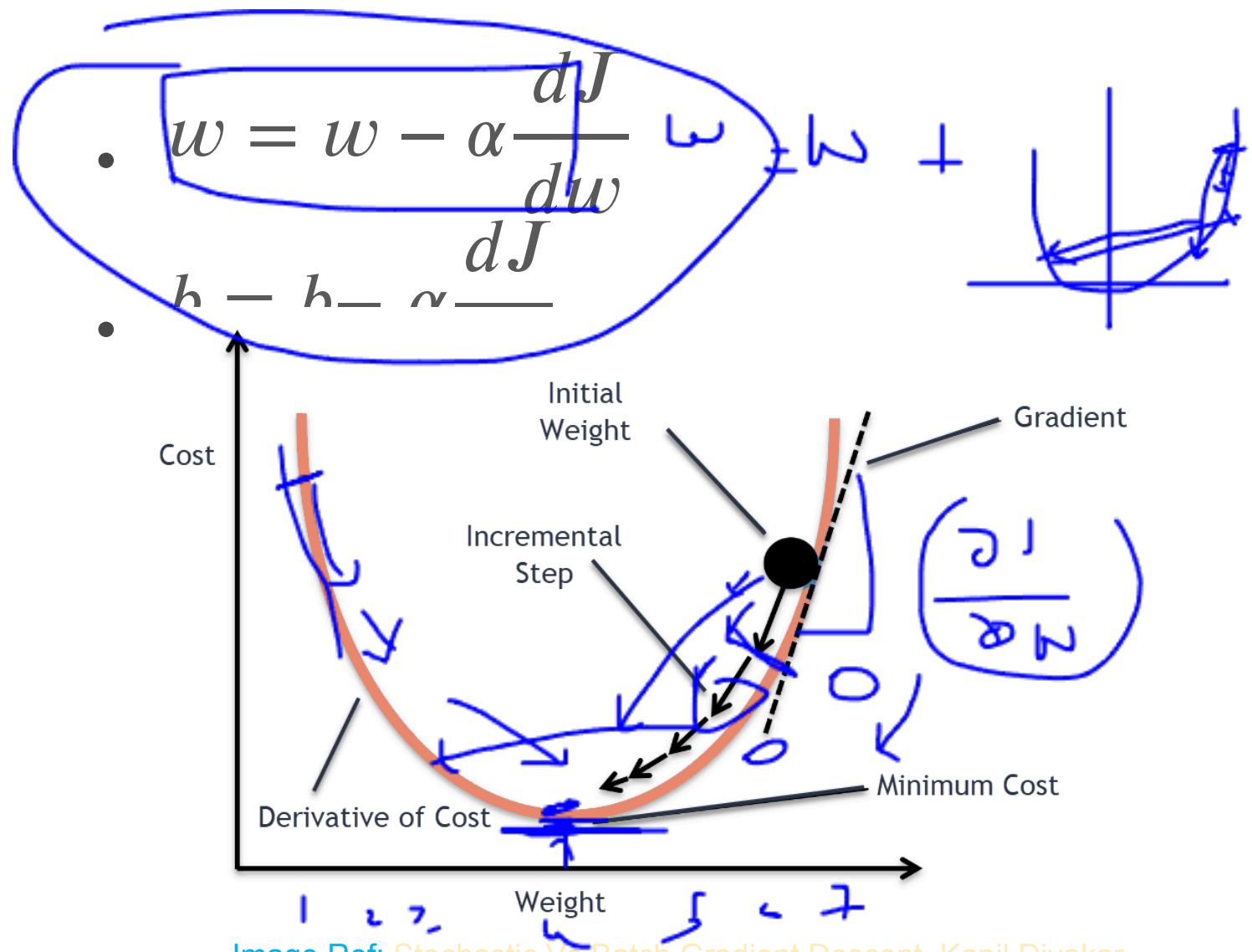


Image Ref: [Stochastic Vs Batch Gradient Descent, Kapil Divakar](#)

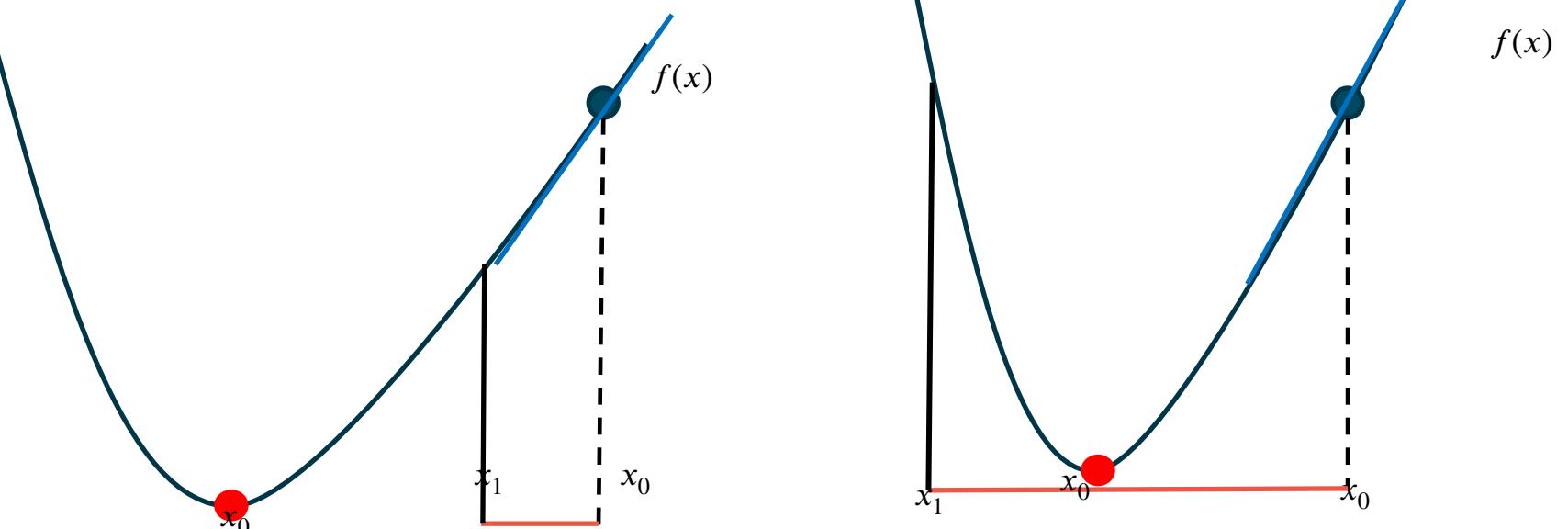
$$\frac{\partial J}{\partial w}, z = \frac{wx + b}{w}$$

$$\frac{\partial J}{\partial w} \rightarrow \text{Partial}$$

$$w = w - \alpha \frac{dJ}{dw}$$

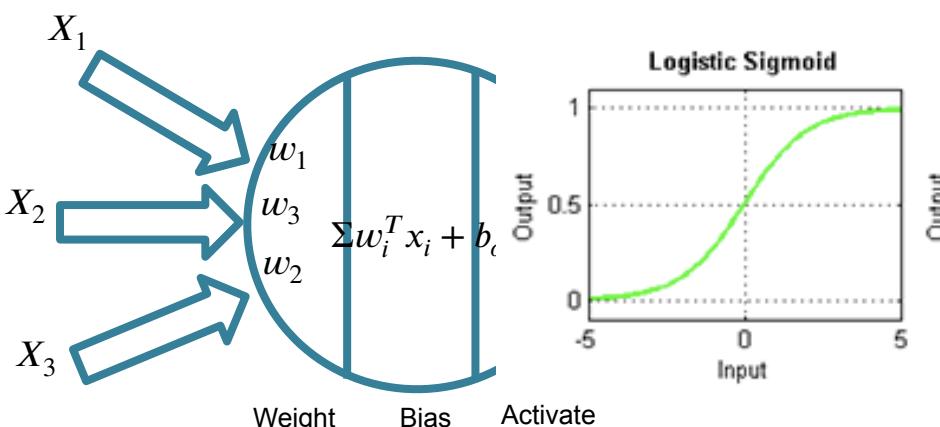
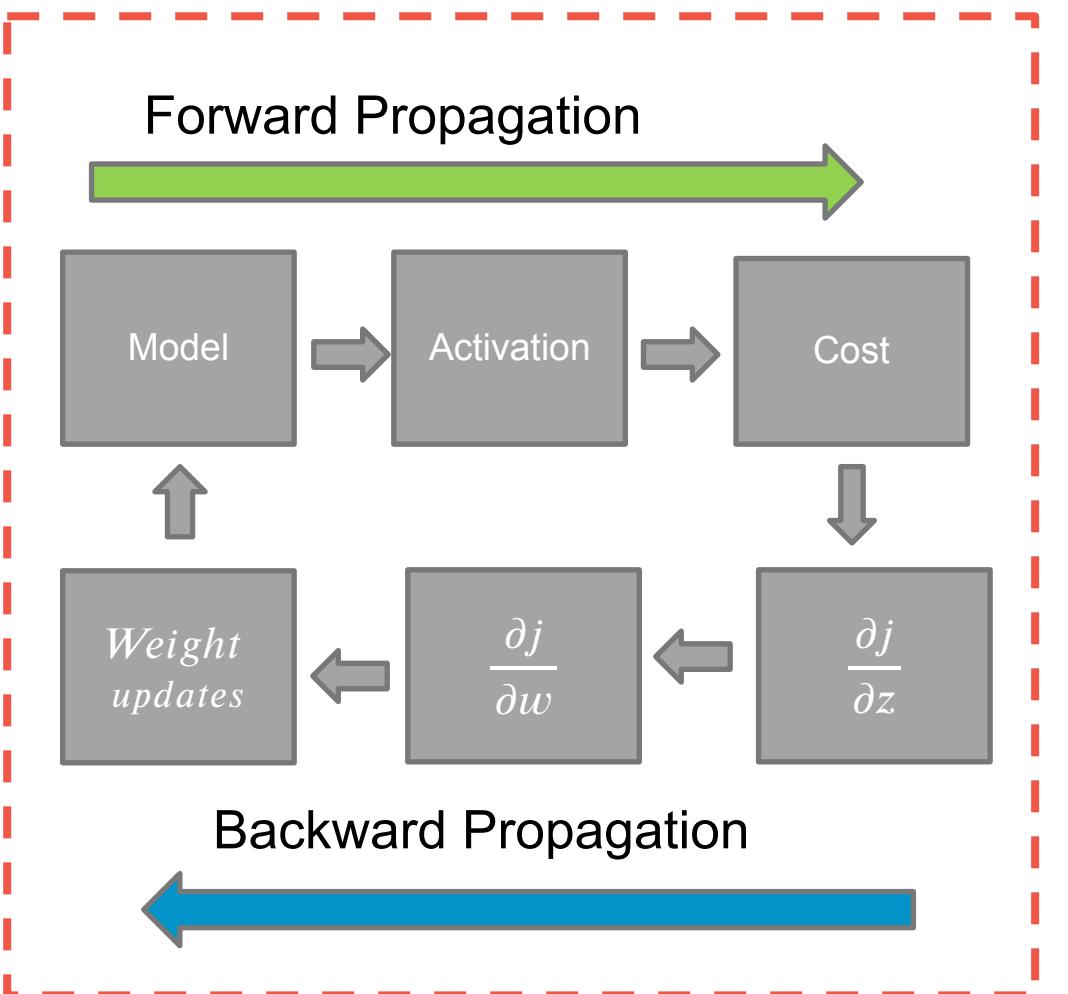
$$\frac{dJ}{dw}$$

- Learning rate α , is between 0 and 1
- The choice of learning rate can be critical to how soon the algorithm converges.



Excel

Forward & Backward Propagation



Cost function (over entire data set)

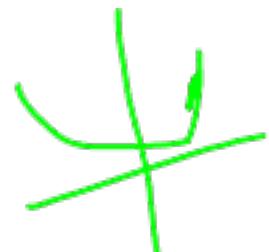
$$J(w, b) = \frac{1}{m} \sum [y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})]$$

Weight update equations

$$\rightarrow w = w - \alpha \frac{dJ}{dw}$$

$$\rightarrow b = b - \alpha \frac{dJ}{db}$$

Learning rate α , is between 0 and 1



→ We need derivatives for weight updates. (Gradient descent equations)

→ Derivates of cost with respect to weights is not easily obtained (analytically or numerically)

→ Computational graph along with chain rule allows you to take derivatives at each step backwards, and then multiply all intermediate derivates.

$$\mathcal{P} = f(w)$$

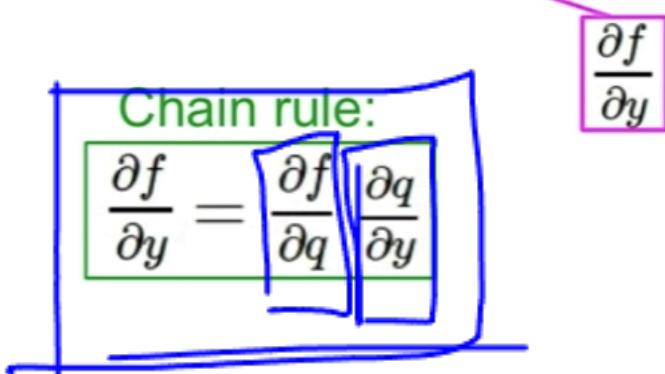
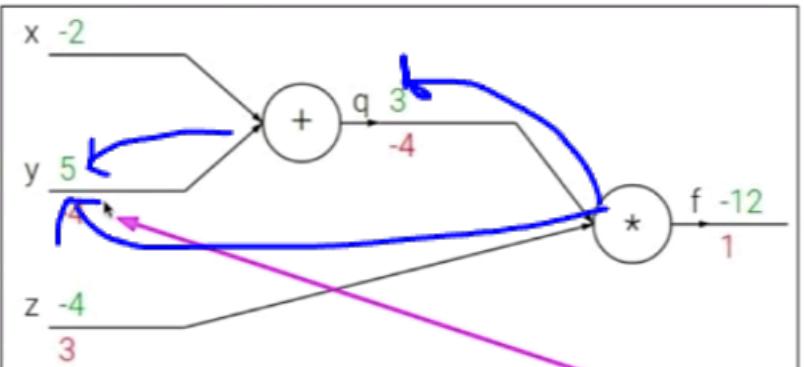
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

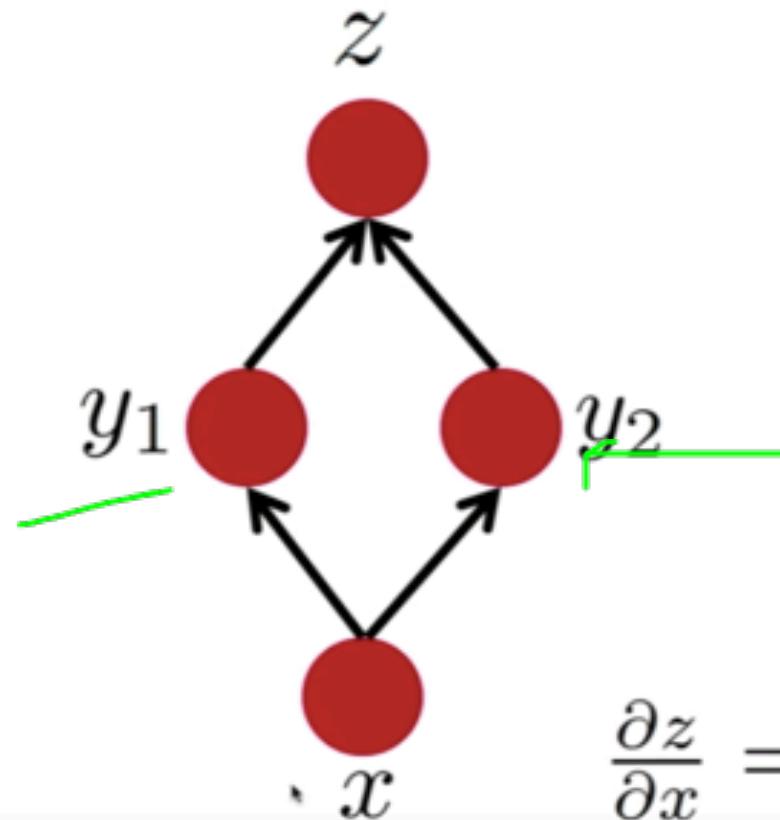


$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$$

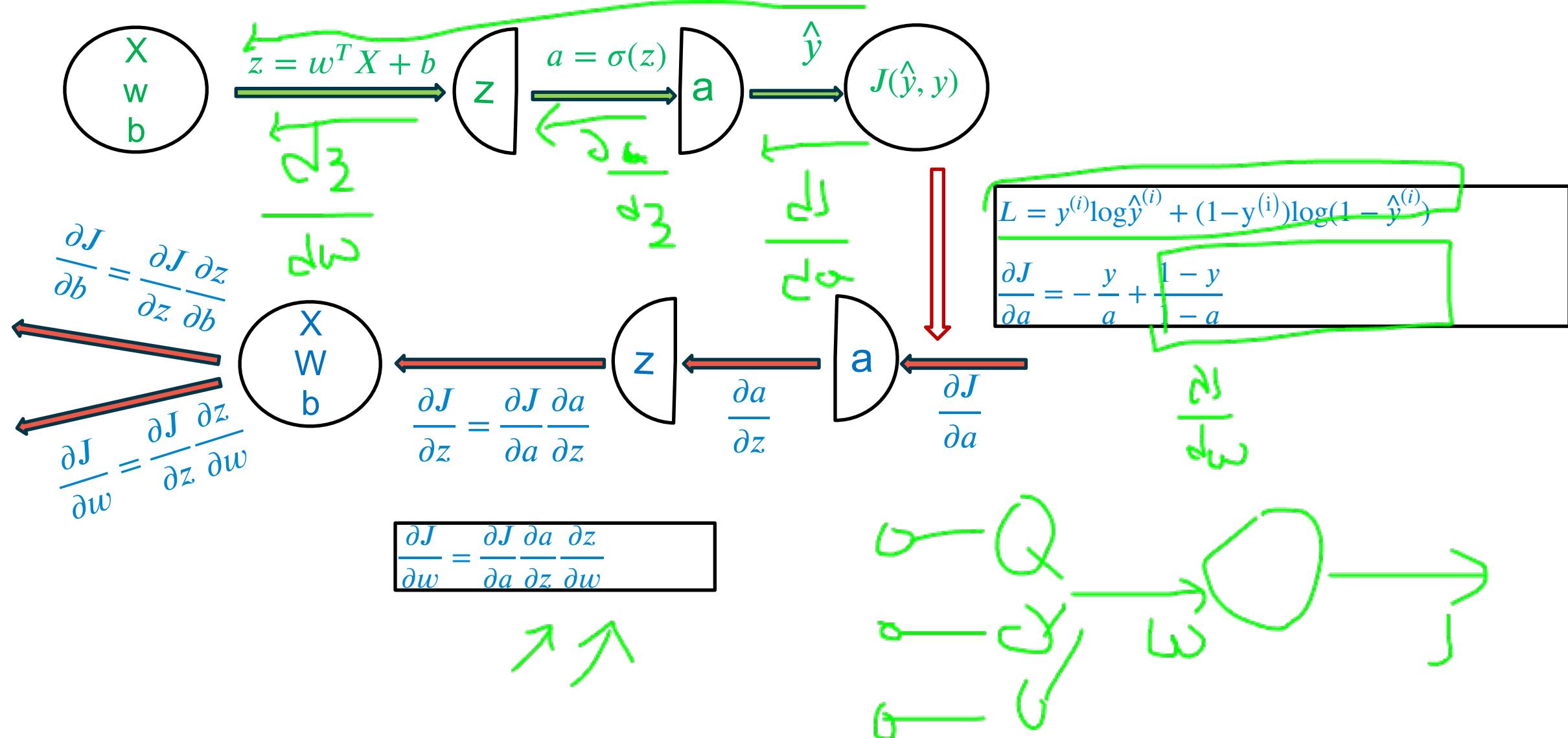
$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} = \frac{\partial f}{\partial y}$$

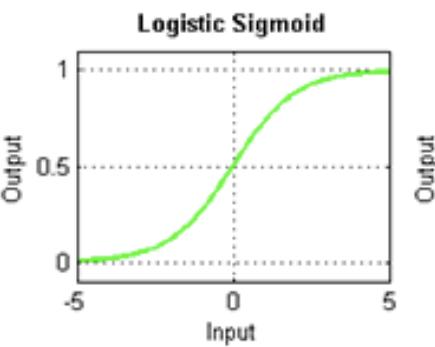
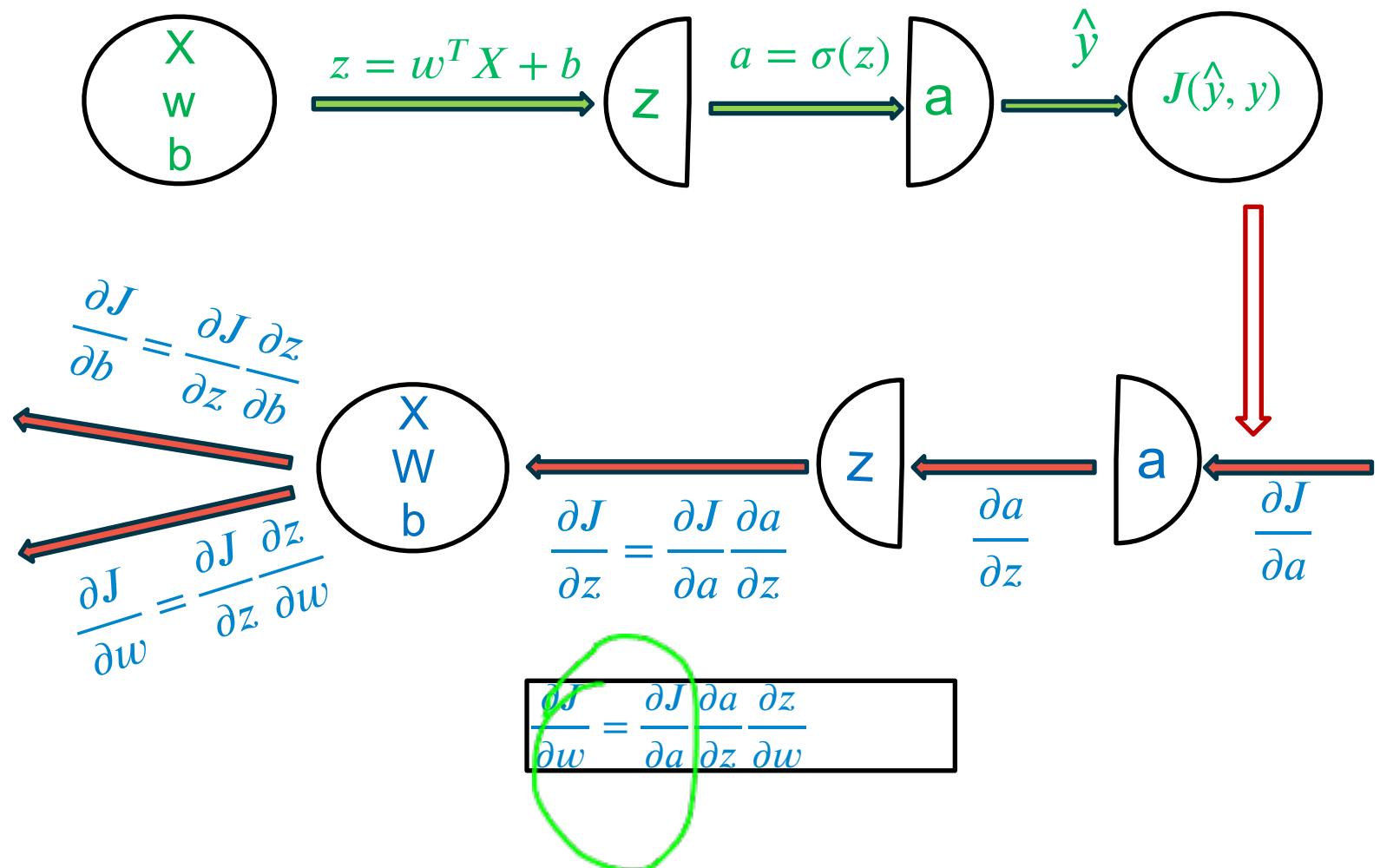
Ref: [Backpropagation and Project Advice](#), Stanford University

Multiple Paths Chain Rule



$$\frac{\partial z}{\partial x} = \underline{\frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial x}} + \boxed{\underline{\frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial x}}}$$





$$a = \frac{1}{1 + e^{-z}}$$

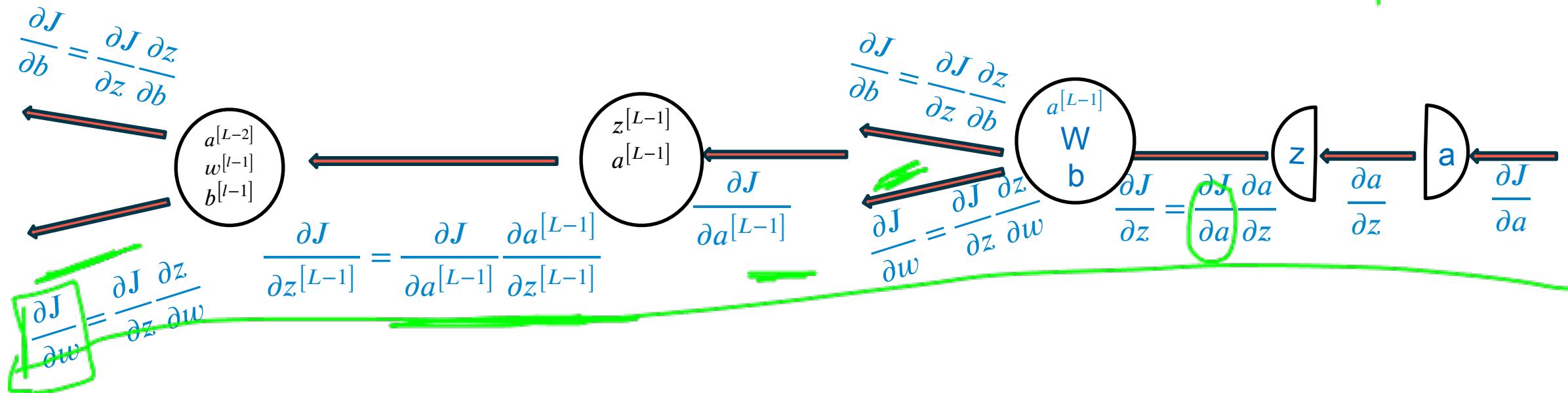
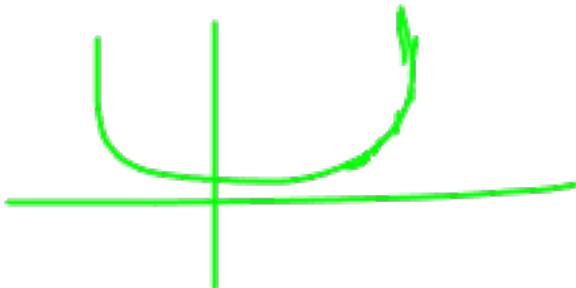
$$\frac{\partial a}{\partial z} = \frac{1}{1 + e^{-z}}(1 - \frac{1}{1 + e^{-z}})$$

$$= a(1 - a)$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w}$$

O

O



Three steps in backpropagation

1. Cost to Activation
2. Activation to Z (Linear combination and bias)
3. Z to weights

$$\frac{\partial J}{\partial a^{[L-1]}} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial a^{[l-1]}} = \frac{\partial J}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial a^{[l-1]}}$$

$$\begin{aligned}
 \frac{\partial J}{\partial w^{[L-1]}} &= \frac{\partial J}{\partial z^{[L-1]}} \frac{\partial z^{[L-1]}}{\partial w^{[L-1]}} \\
 &= \frac{\partial J}{\partial a^{[L-1]}} \frac{\partial a^{[L-1]}}{\partial z^{[L-1]}} \frac{\partial z^{[L-1]}}{\partial w^{[L-1]}} \\
 &= \frac{\partial J}{\partial z^{[l-1]}} \frac{\partial z^{[l-1]}}{\partial a^{[l-1]}} \frac{\partial a^{[L-1]}}{\partial z^{[L-1]}} \frac{\partial z^{[L-1]}}{\partial w^{[L-1]}} \\
 &\quad \text{---} \\
 &= \frac{\partial J}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial a^{[l-1]}} \frac{\partial a^{[L-1]}}{\partial z^{[L-1]}} \frac{\partial z^{[L-1]}}{\partial w^{[L-1]}}
 \end{aligned}$$

$$1. \frac{\partial J}{\partial z} = \frac{\partial J}{\partial a} \frac{\partial a}{\partial z} = a - y \quad (\text{see note 1})$$

$$2. \frac{\partial J}{\partial w} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial w} = \frac{\partial J}{\partial z} x^T. \quad (\text{see note 2})$$

$$3. \frac{\partial J}{\partial b} = \frac{\partial J}{\partial z}$$

$$4. \frac{\partial J}{\partial z^{[L-1]}} = w \frac{\partial J}{\partial z} * \frac{\partial a}{\partial z}$$

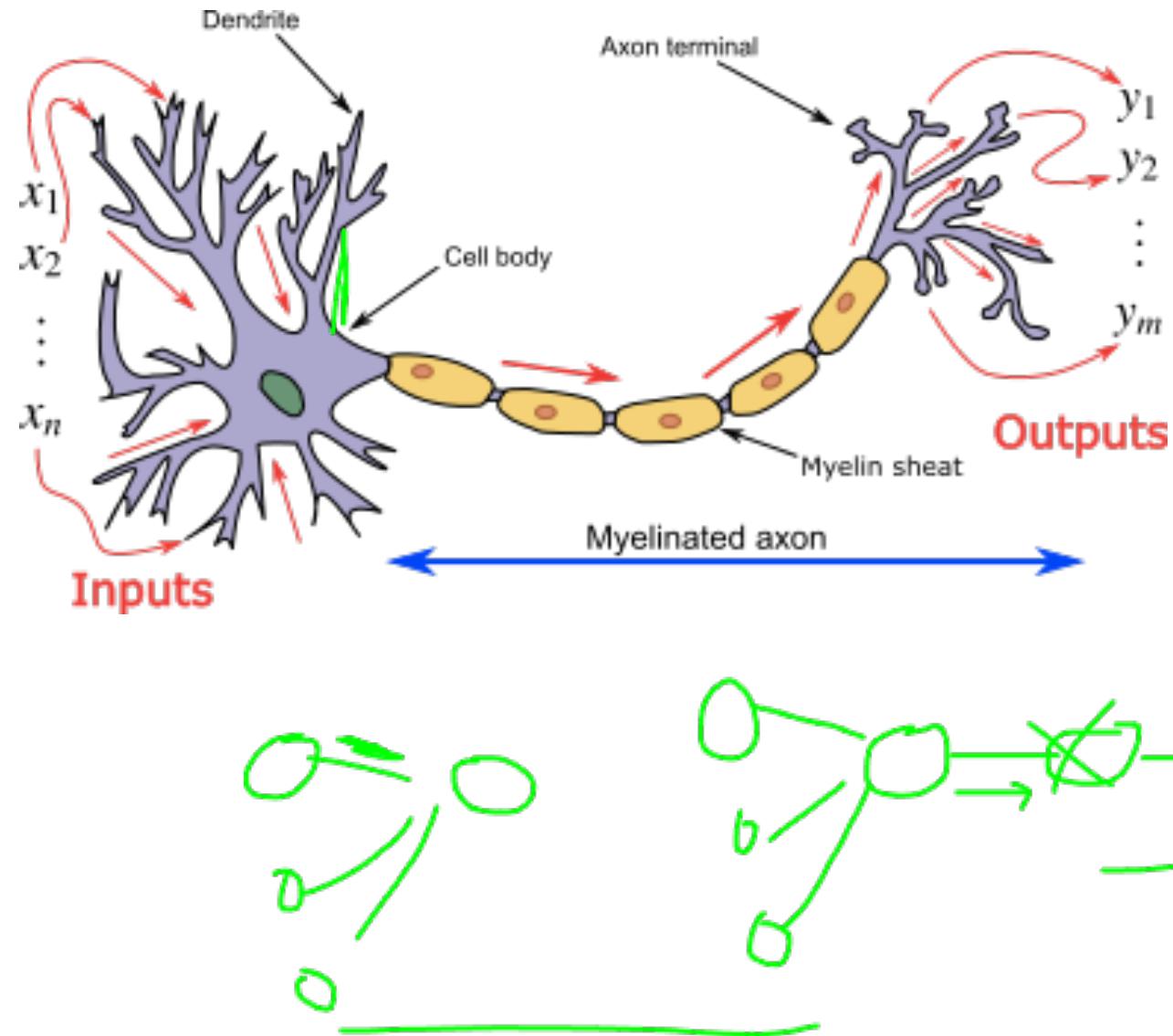
$$5. \frac{\partial J}{\partial w^{[L-1]}} = \frac{\partial J}{\partial z^{[L-1]}} \frac{\partial z^{[L-1]}}{\partial w^{[L-1]}} = \frac{\partial J^{[L-1]}}{\partial z^{[L-1]}} x^T$$

$$6. \frac{\partial J}{\partial b^{[L-1]}} = \frac{\partial J}{\partial z^{[L-1]}}$$

(Note 1: Only for sigmoid. Will differ for other activation functions)

(Note 2: If this not the input layer, use (a) of previous layer, instead of x)

- Each step of gradient descent requires gradient to be calculated over the entire data
- Sending ALL the data at once (if in millions), could be computationally infeasible
- Alternate is to select a ‘random’ data point, calculate the derivative and update the weight. Then take another point and repeat the process.
- Since only one observation is taken, the gradient descent may fluctuate a lot and may take more iterations (but less computations) to converge.
- Sometimes, instead of one data point, a “**mini-batch**” of data points is selected which is faster than stochastic gradient descent and at the same time less computation than “batch” (full) gradient descent
- When the entire data set is passed through the network once, it is called an **Epoch**,



Goal Understand the brain

Emulate the Brain

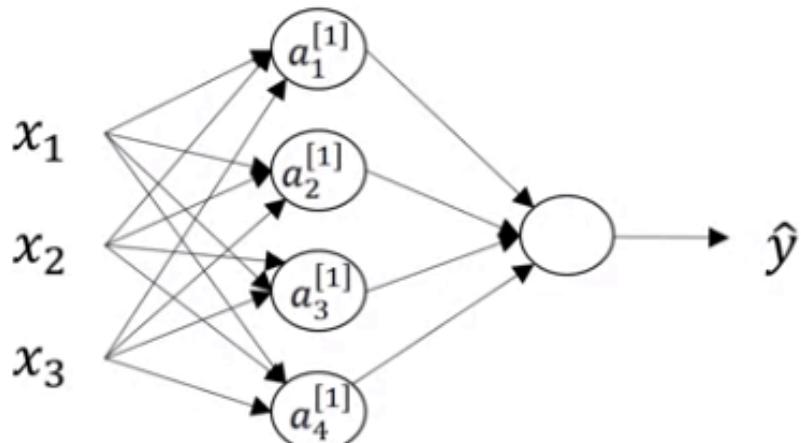
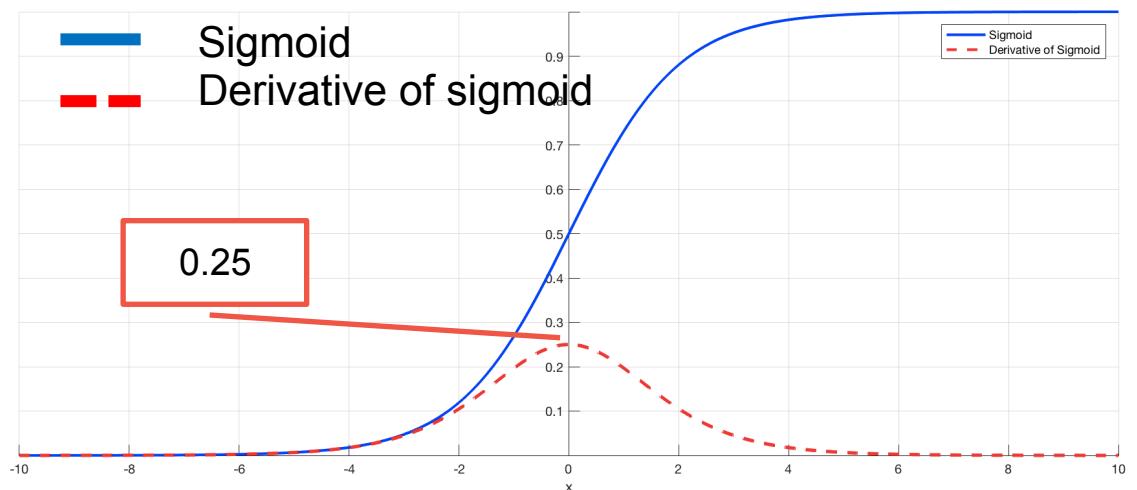
ANN Did we achieve the goal?

Exploding & Vanishing Gradient

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w}$$

$$\frac{\partial J}{\partial w^{[L-1]}} = \frac{\partial J}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial a^{[l-1]}} \frac{\partial a^{[L-1]}}{\partial z^{[L-1]}} \frac{\partial z^{[L-1]}}{\partial w^{[L-1]}}$$

- Deeper you go in the network, weights of initial layers will have several gradients multiplied (due to backpropagation)
- Gradient of sigmoid activation function with max value of 0.25, can vanish (and hence not update weights) in a deeper network.
- Typical problem in sigmoid but not necessarily in RELU



- In some cases gradients may explode or vanish if weights are not initialized properly. This is not because of sigmoid function but similar reasoning apply (i.e. several weight matrices are multiplied over the layers in a deep network)

$$\frac{\partial J}{\partial z^{[L-1]}} = w \frac{\partial J}{\partial z} * \frac{\partial a}{\partial z}$$

$$\frac{\partial J}{\partial w^{[L-1]}} = \frac{\partial J}{\partial z^{[L-1]}} \frac{\partial z^{[L-1]}}{\partial w^{[L-1]}} = \frac{\partial J}{\partial z^{[L-1]}}$$

$$\frac{\partial J}{\partial b^{[L-1]}} = \frac{\partial J}{\partial z^{[L-1]}}$$

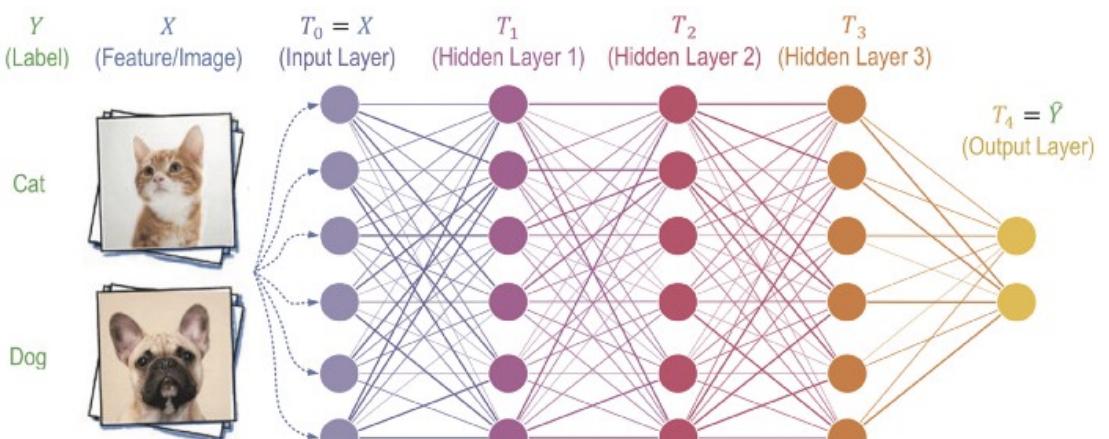
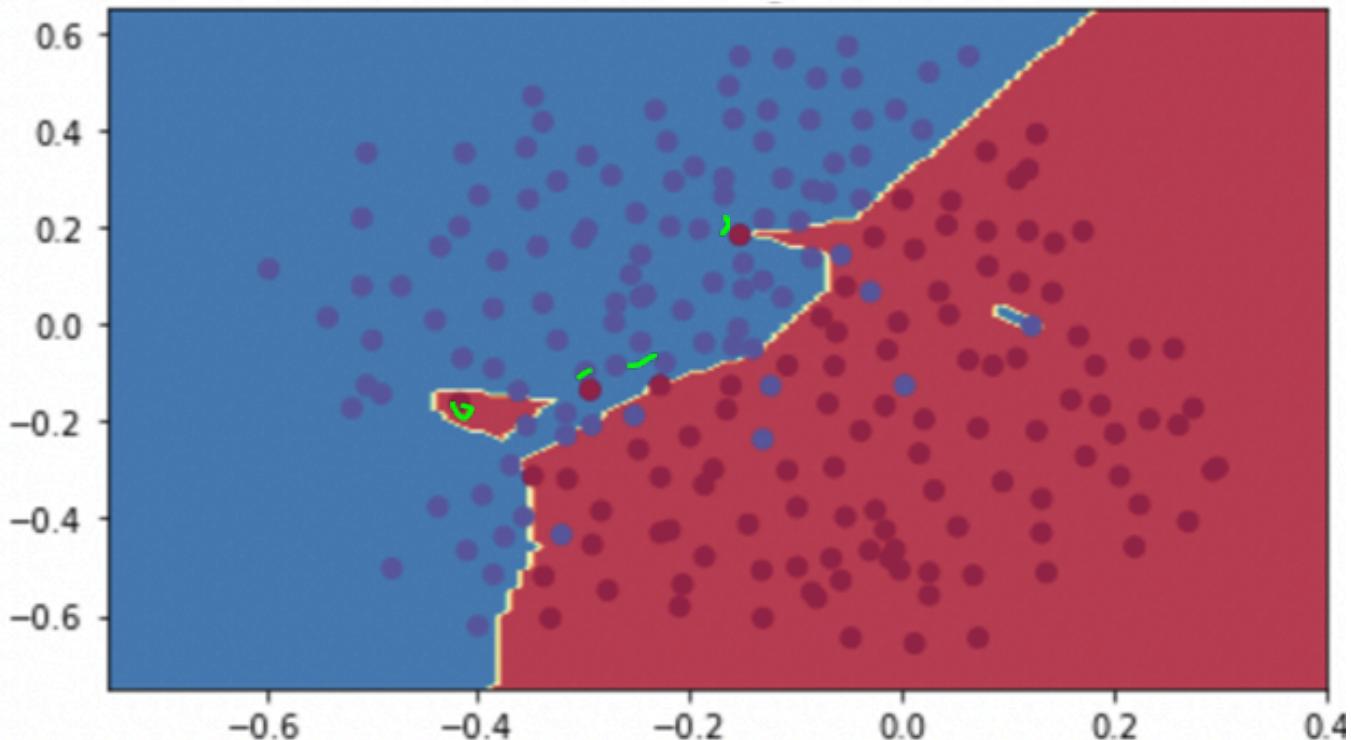


Image [source](#)

Regularization

- Neural Networks can overfit!

Model without regularization



$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

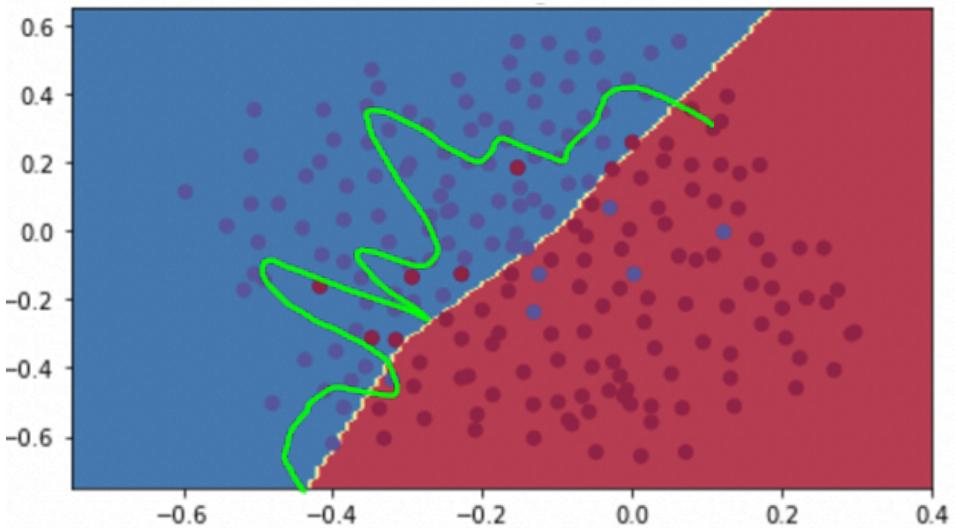
L2 regularization

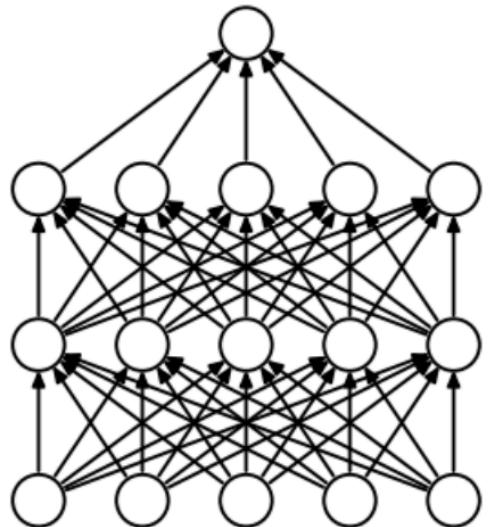
Regularization Parameter

Image [source](#)

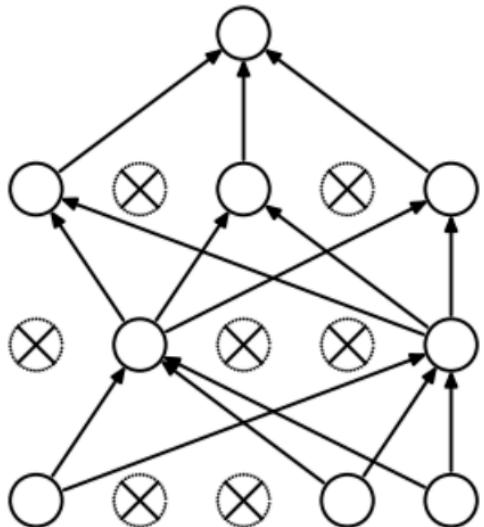
As the weights increase (i.e high weight is given to any neuron, regularization term also increases and adds to the cost. Since, gradient descent tries to minimize the cost, it will try to keep the weights low, hence regularized)

Model with regularization





(a) Standard Neural Net

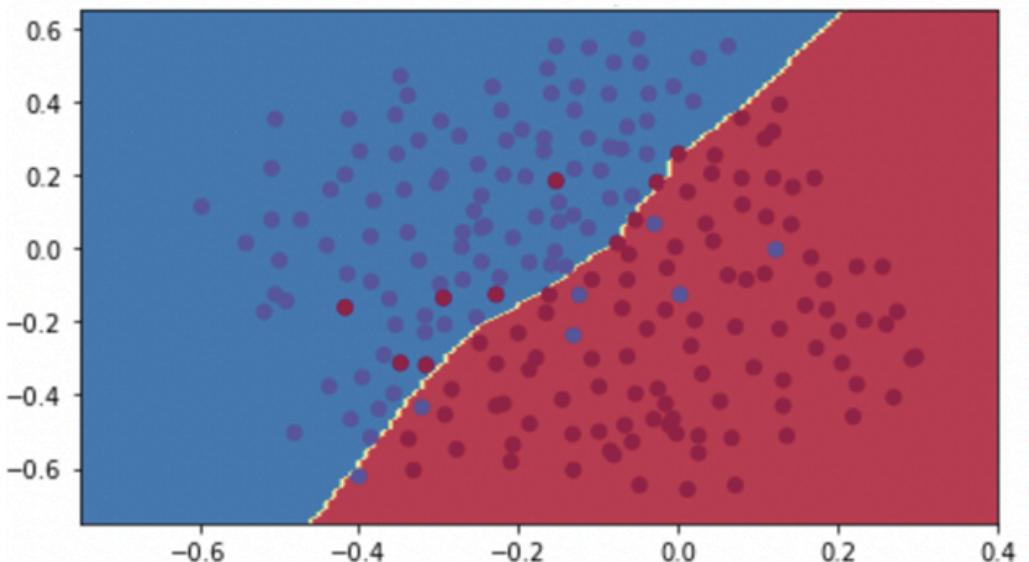


(b) After applying dropout.

Dropout Regularization

Random neurons are dropped from the network with some probability. Since a neuron has certain chances of being dropped, network will not give high weight to any of it during training providing the effect of regularization.

Model with dropout



➤ **Epoch:**

- One Epoch is when an ENTIRE dataset is passed forward and backward through the neural network only ONCE.

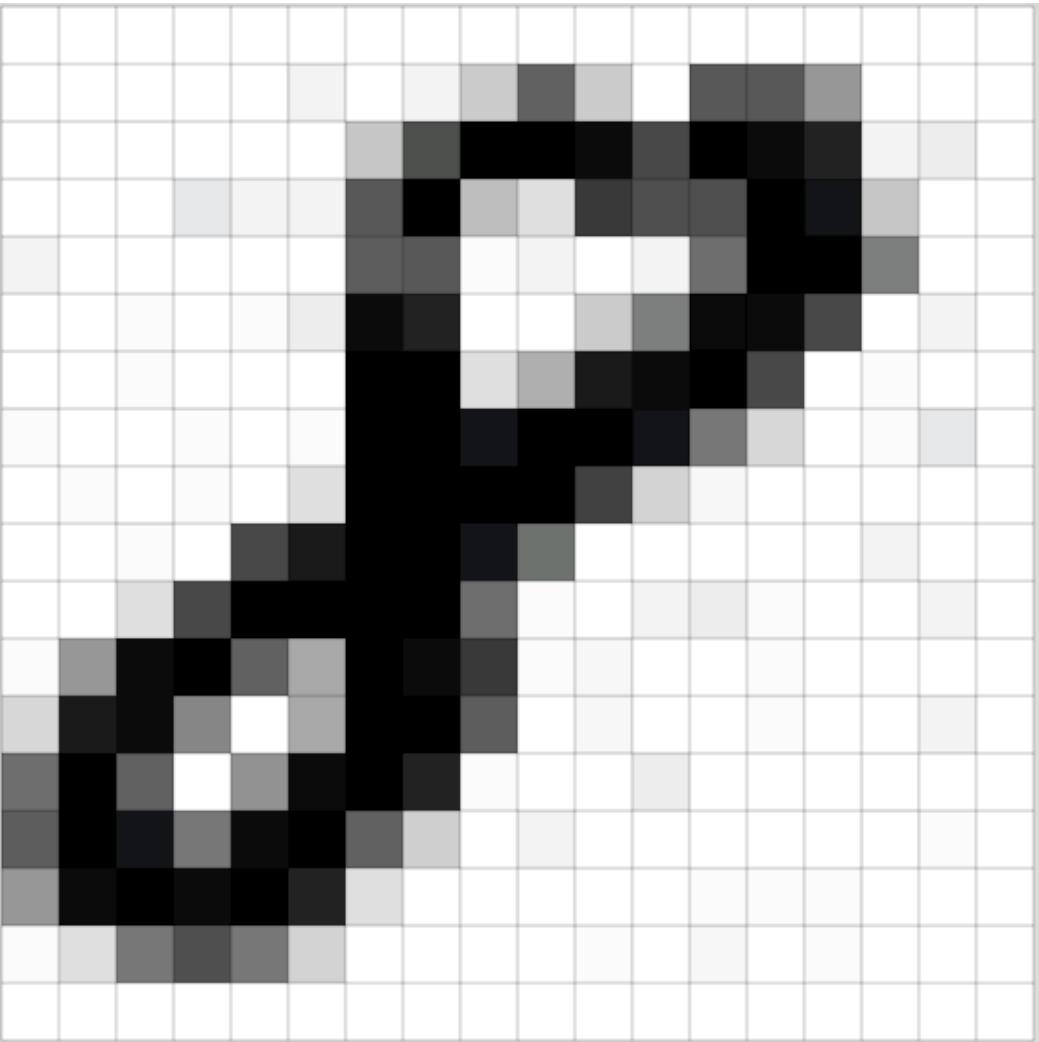
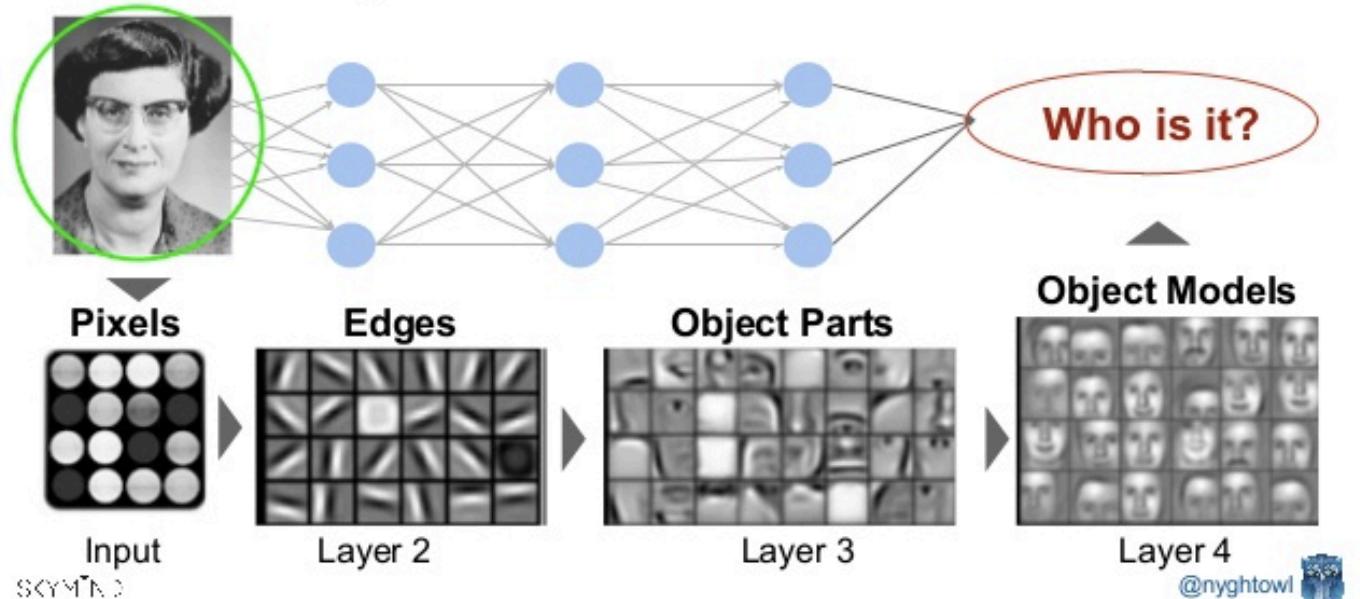
➤ **Batch:**

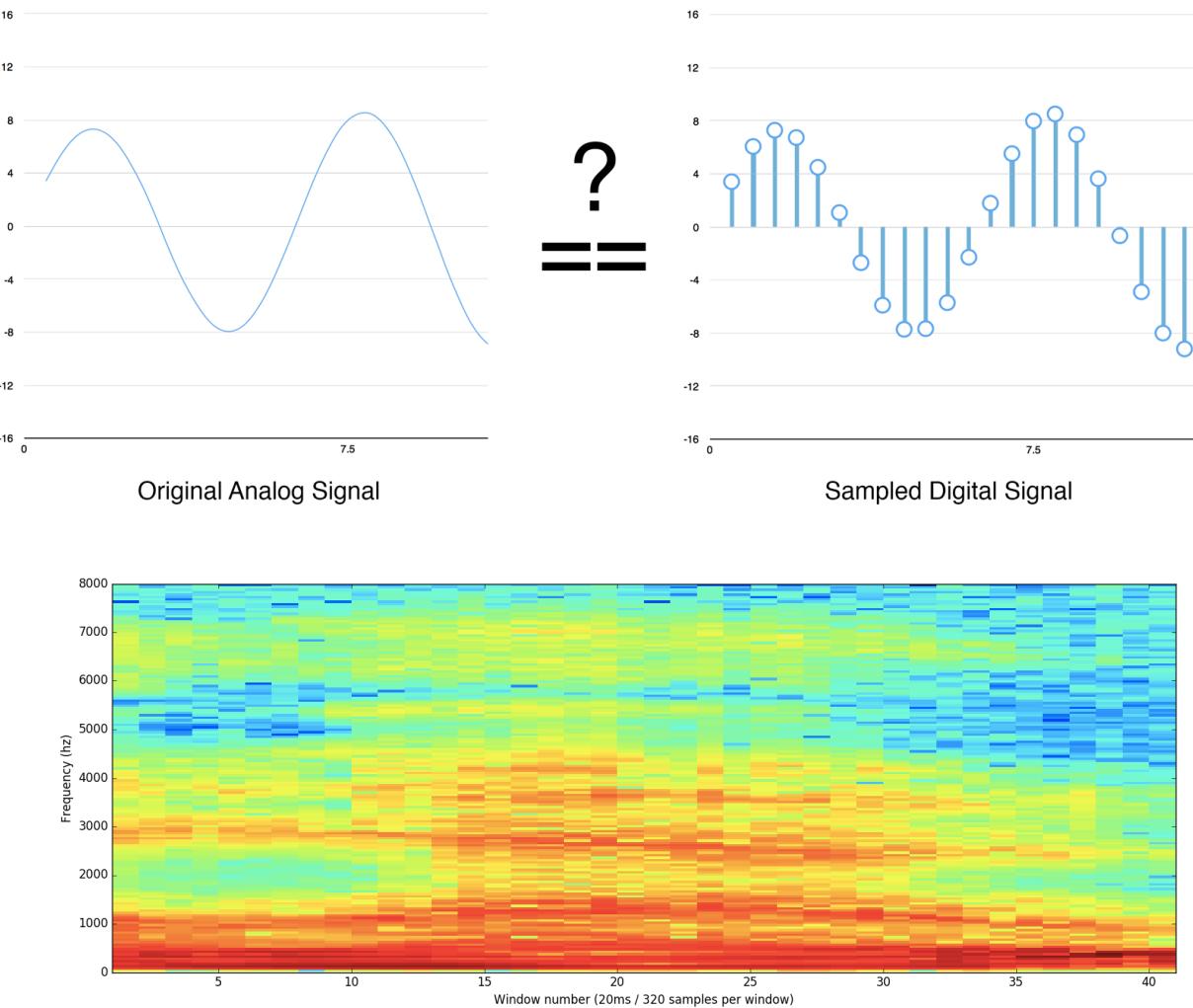
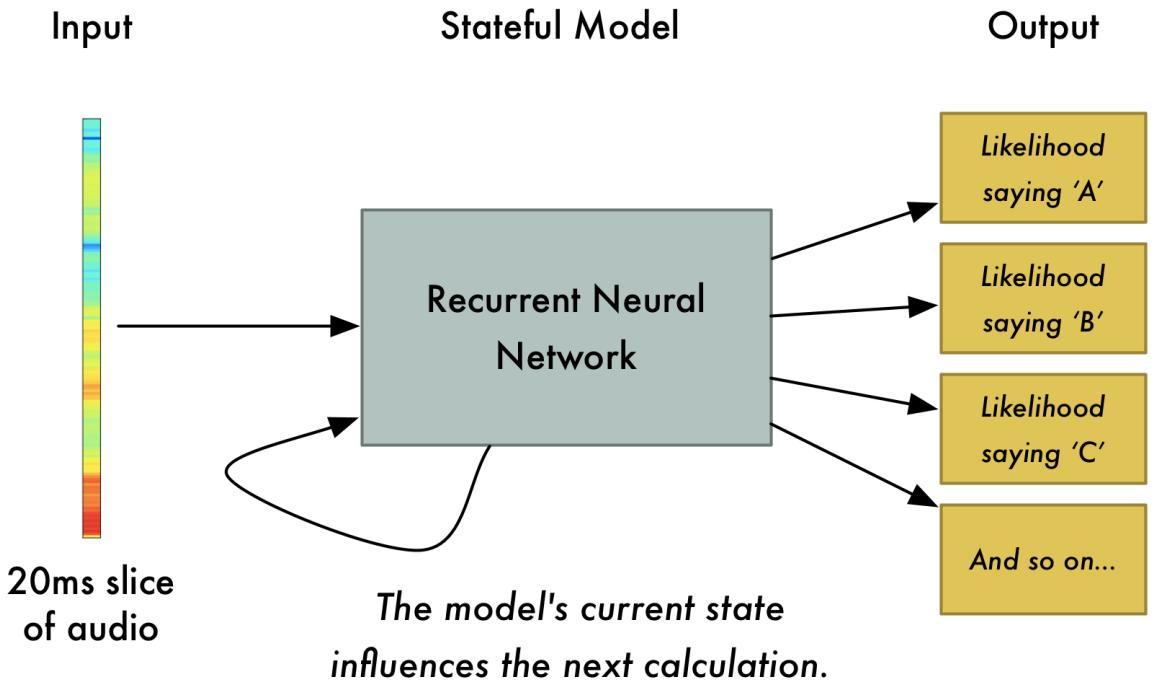
- Total number of training examples present in a single batch.
- If you can't pass the entire dataset into the neural net at once, divide dataset into Number of Batches or sets or parts.

- **Batch Gradient Descent.** Batch Size = Size of Training Set
- **Stochastic Gradient Descent.** Batch Size = 1
- **Mini-Batch Gradient Descent.** $1 < \text{Batch Size} < \text{Size of Training Set}$

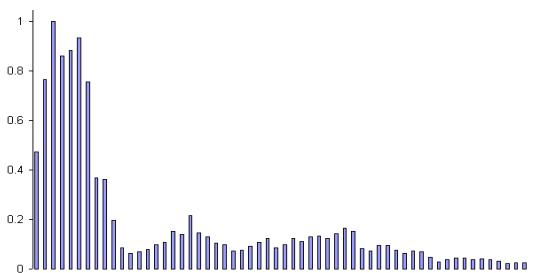
Use Cases

NN Example





- Task: Learn to discriminate between two different voices saying “Hello”
- Data
 - Sources
 - Steve Simpson
 - David Raubenheimer
 - Format
 - Frequency distribution (60 bins)
 - Analogy: cochlea 1sw 1sw



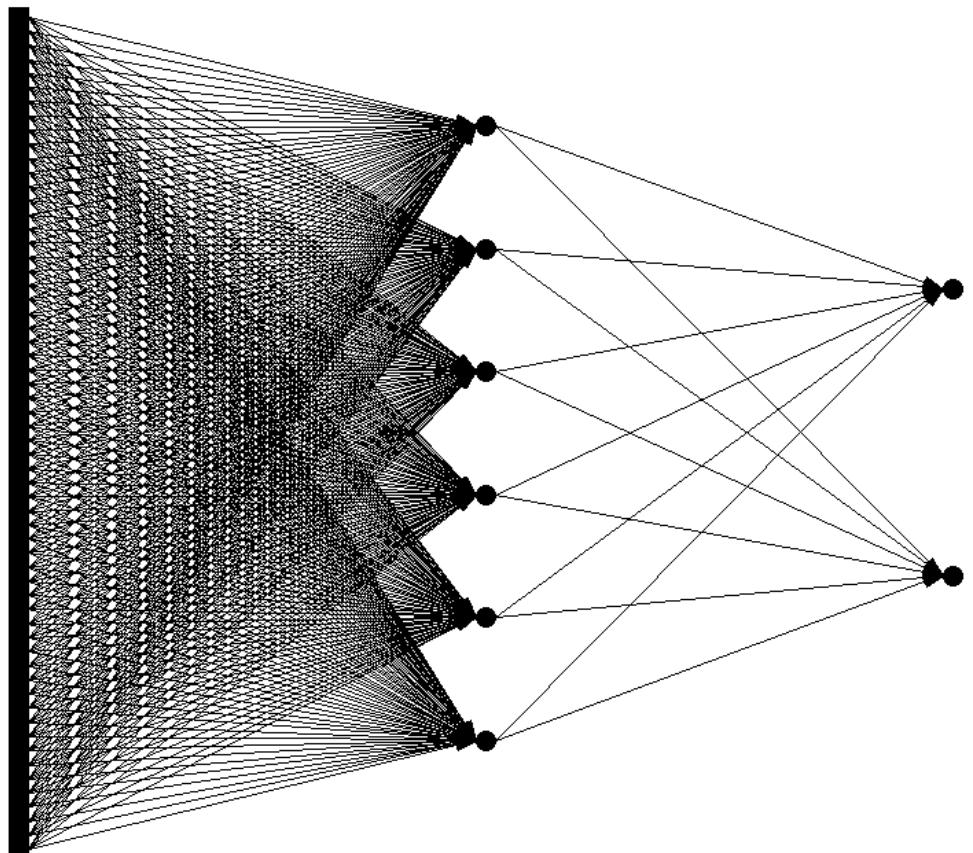
Network architecture

Feed forward network

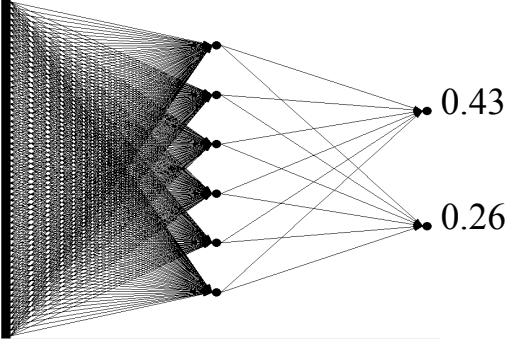
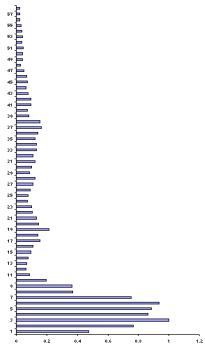
60 input (one for each frequency bin)

6 hidden

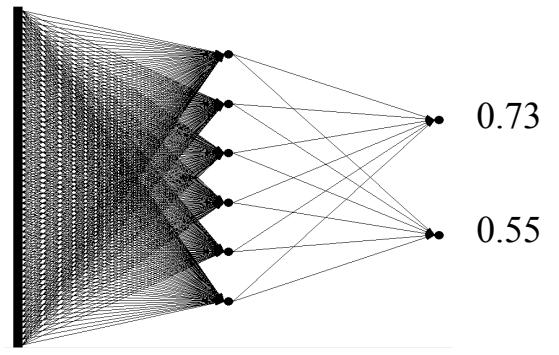
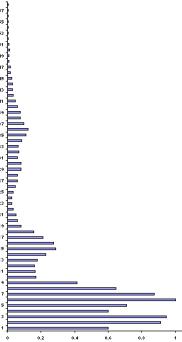
2 output (0-1 for “Steve”, 1-0 for “David”)



Steve

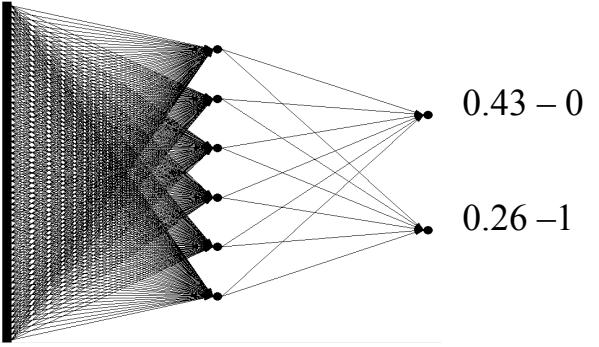
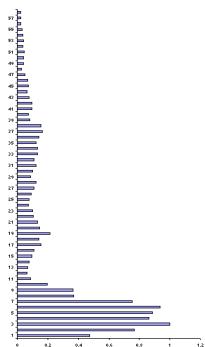


David



Calculate error

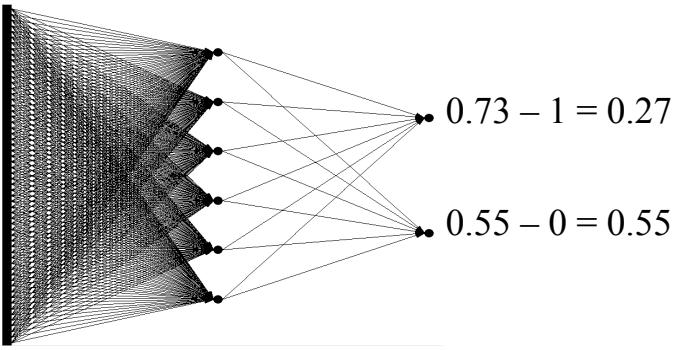
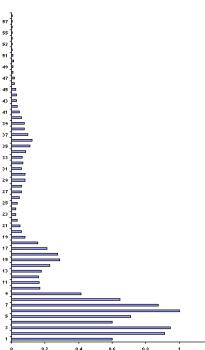
Steve



David

$$= 0.43$$

$$= 0.74$$

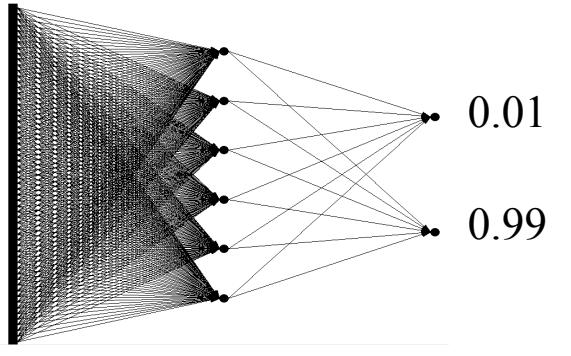
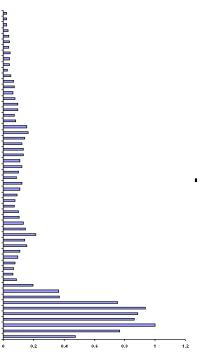


Performance of trained network

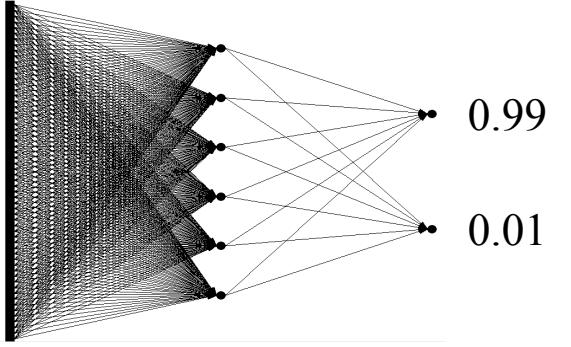
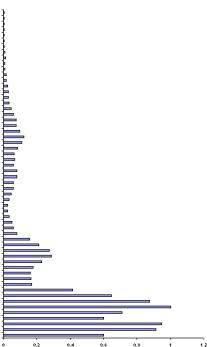
Discrimination accuracy
between known “Hello”s
100%

Discrimination accuracy
between new “Hello”s
100%

Steve



David



Back Up