

## Topics For Today

1. Recap of multiclass classification

2. Deep Neural Networks (Back prop)

- Scalar chain }
- Back prop for a MLP → Multi Layer Perceptron
- Activation functions

• Demo

• Dropout

3. Variants of Gradient Descent + Demo

"Vanilla" Gradient Descent

# Multiclass classification

$n \rightarrow$  no. of features  $\rightarrow$

$m \rightarrow$  No of examples / data pts



$K = 2$  Binary Classification



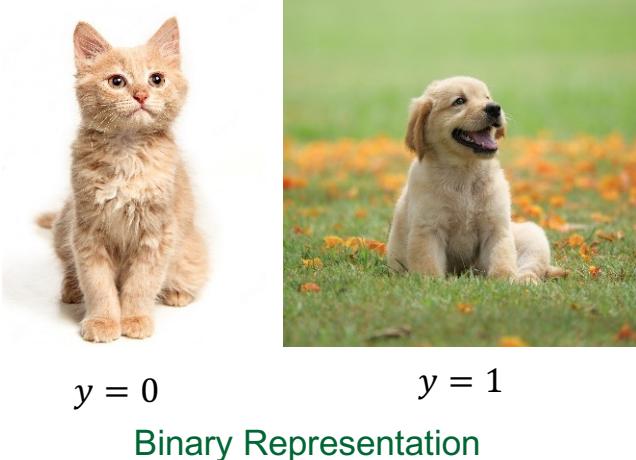
$K = 3$  Multiclass Classification

- $K$  : Number of classes.
  - $K = 2$  Binary classification
  - $K > 2$  Multiclass (Multinomial) classification
- **ML Recipe** – Requires the following
  - **Data Representation** – How do we represent the output (it cannot be 0 or 1)
  - **Model** – What is the nonlinearity instead of sigmoid? (That was used to lie between 0 and 1)
  - **Loss** – What is the Loss function?
  - **Optimization** – Gradient Descent but What is the gradient for the new model?

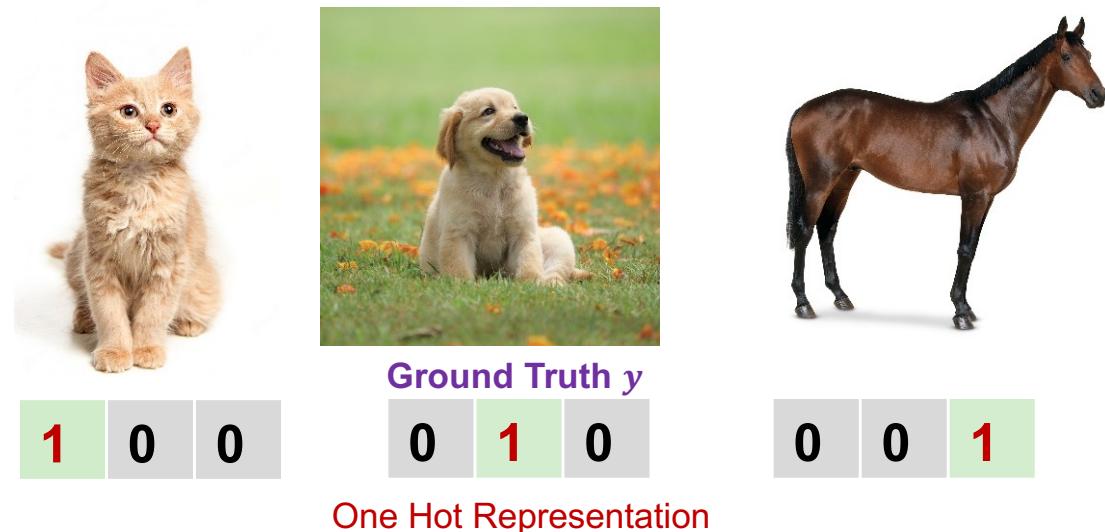
Let us look at these issues one by one

# Data Representation – One Hot Vector

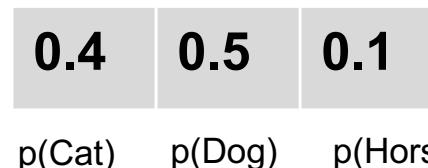
$K = 2$  Binary Classification



$K = 3$  Multiclass Classification



- “One-Hot Representation” – Represent the output as a vector instead of a scalar
- Each output class is represented by a 1 in one location and 0 in others – One Hot
- We can interpret the entries in each box as the probability of the class



Model Output  $\hat{y}$   
representation

$$\sum_{i=1}^3 \hat{y}_i = 1$$

For both  $y$  and  $\hat{y}$ , every entry represents  $p(y_i = 1|x)$

# Model Nonlinearity -- Softmax

$\hat{y}_1$   
 $\hat{y}_2$   
 $\hat{y}_3$



0.4    0.5    0.1

p(Cat)    p(Dog)    p(Horse)

$$\begin{aligned}\hat{y} &= [0.4 \quad 0.5 \quad 0.1] \\ &= [\hat{y}_1 \quad \hat{y}_2 \quad \hat{y}_3]\end{aligned}$$

	$z$	$\sigma(z)$	Softmax(z)
1	1	0.7311	0.1191
2	3	0.9526	0.8801
3	-4	0.0180	0.0008
SUM $\Sigma$		1.7016	1.0000

- We can define  $z = \sum w_i x_i$  as before
- However, we cannot use  $\hat{y} = \sigma(z)$  since we would not have  $\sum \hat{y}_i = 1$
- To guarantee that the output can be interpreted as a probability we introduce softmax.
- This is defined as

$$sf(z_1) = \frac{e^{z_1}}{D}, sf(z_2) = \frac{e^{z_2}}{D} \Rightarrow sf(z_1) + sf(z_2) + sf(z_3) = 1$$

$$softmax(z_1) = \frac{\exp(z_1)}{(exp(z_1) + exp(z_2) + exp(z_3))}$$

- This has all outputs between 0 and 1 as well as  $\sum \hat{y}_i = 1$

# Categorical Cross-Entropy (CCE) Loss Function

- Recall the binary cross entropy loss function

$$J^{BCE}(y, \hat{y}) = -y \ln \hat{y} - (1 - y) \ln (1 - \hat{y})$$

- We use a generalization of this for the multiclass case

$$J^{CCE}(y, \hat{y}) = -\sum_{i=1}^K (y_i \ln \hat{y}_i)$$

$$-y_1 \ln \hat{y}_1 - y_2 \ln \hat{y}_2 - y_3 \ln \hat{y}_3 = -1 * \ln(0.2)$$

- What happens if apply CCE to the binary case with one-hot?

- $J^{CCE}(y, \hat{y}) = -\sum_{i=1}^K (y_i \ln \hat{y}_i) = -y_1 \ln \hat{y}_1 - y_2 \ln \hat{y}_2$
- In the binary case  $\hat{y}_2 = (1 - \hat{y}_1)$  and  $y_2 = (1 - y_1)$   
 $\Rightarrow J^{CCE}(y, \hat{y}) = -y_1 \ln \hat{y}_1 - (1 - y_1) \ln(1 - \hat{y}_1)$
- This is identical to the binary cross entropy loss!

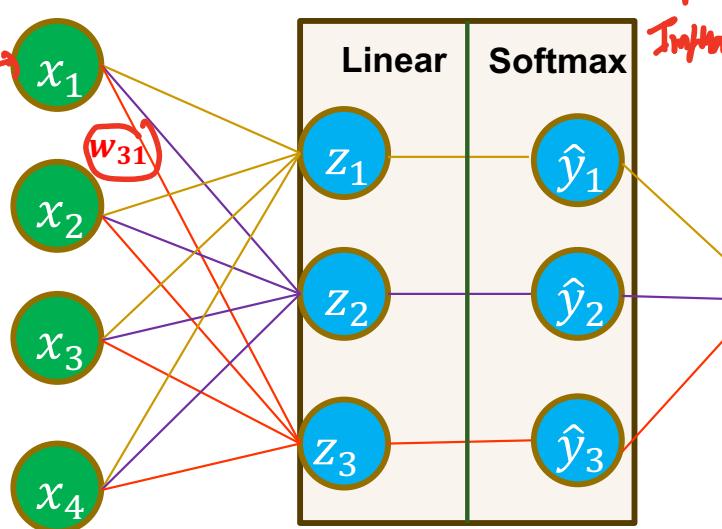


1	Binary	0
1   0	One-hot	0   1

# Optimization – Gradients in the Multiclass case

1.

2.

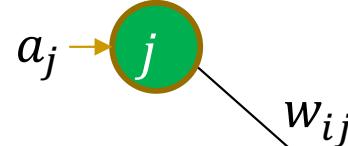


$$w_{31} = \cancel{w_{31}} \quad \alpha$$

↑  
Input  
↓  
Output

$$\frac{\partial J}{\partial w_{31}}$$

$a_j = x_j$  Activation of Input



$\delta_i = \hat{y}_i - y_i$  Error in output

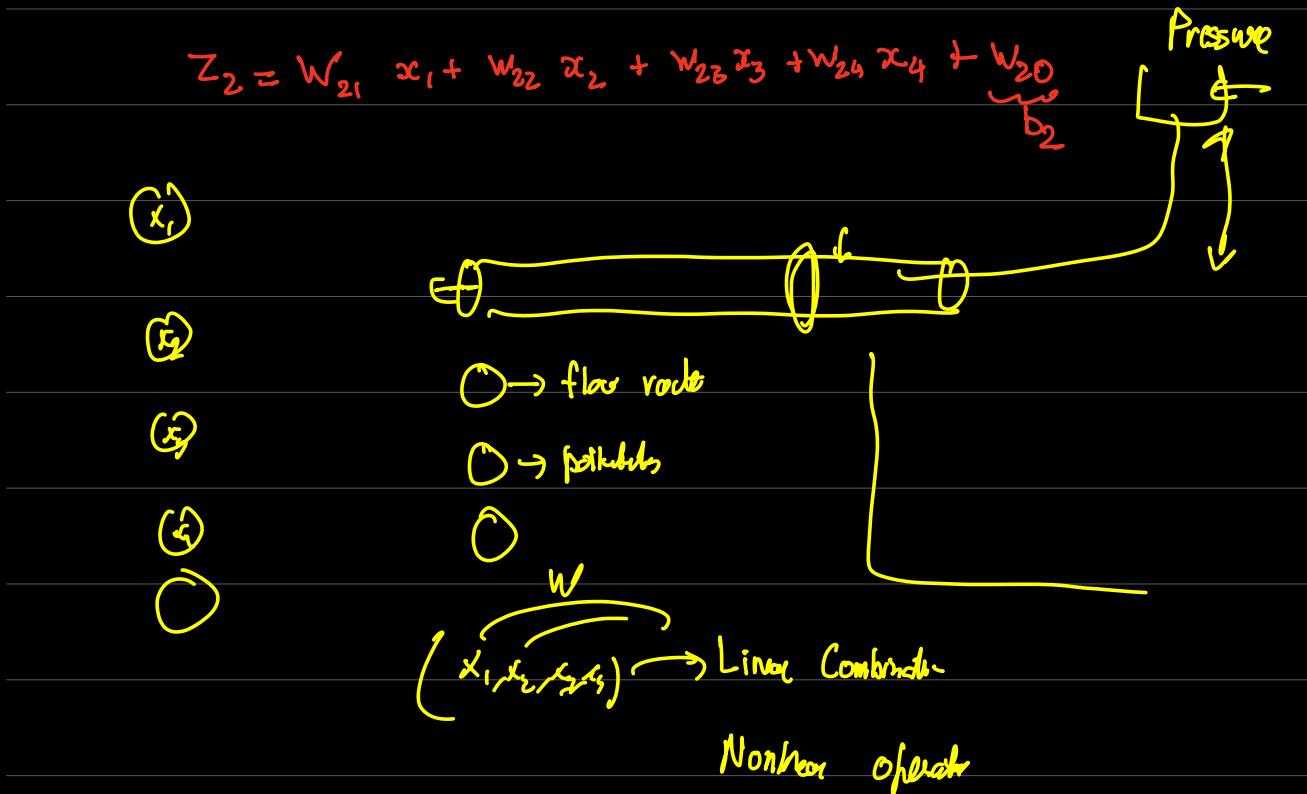
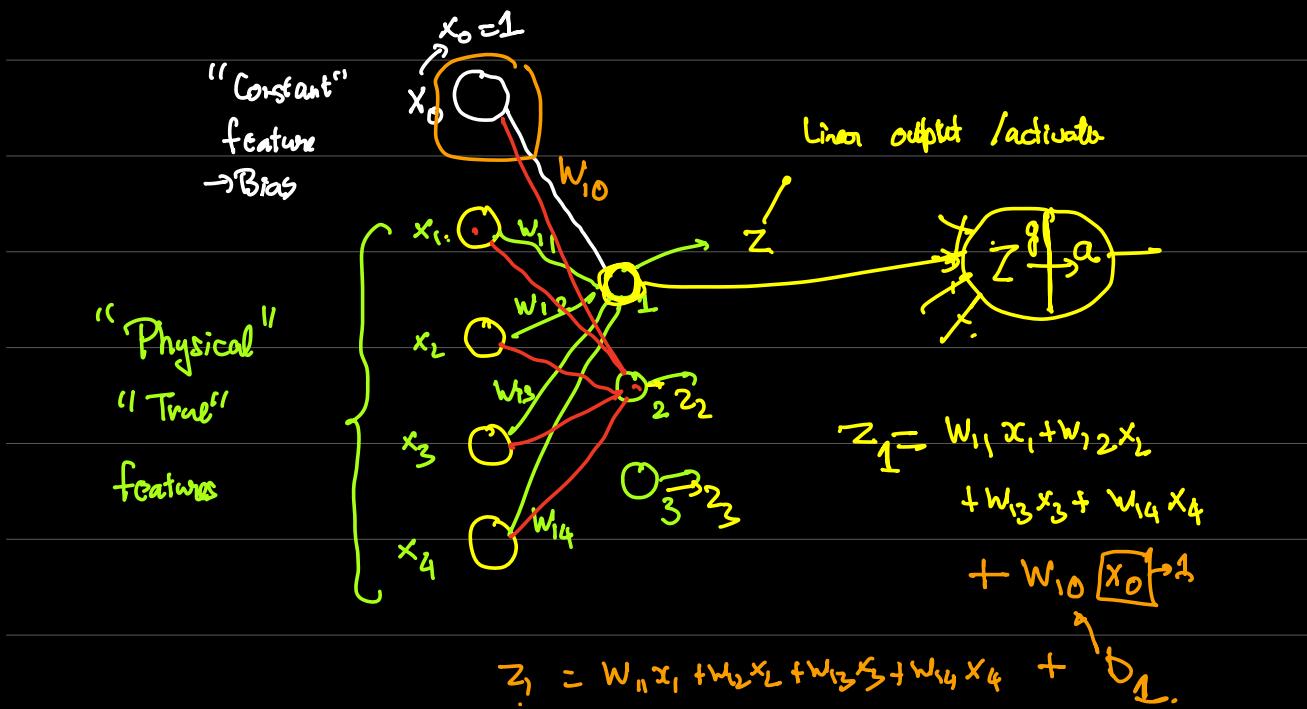
$$(2 - 1)*2$$

- It can be shown that  $\frac{\partial J}{\partial z_1} = \hat{y}_1 - y_1$ . Similarly,  $\frac{\partial J}{\partial z_i} = \hat{y}_i - y_i$
- Same  $\frac{\partial J}{\partial z}$  expression as linear and logistic
- $\frac{\partial J}{\partial w_{31}} = \text{Error in Output } 3 \times \text{Activation of Input } 1 = \delta_3 \times a_1 = (\hat{y}_3 - y_3) \times x_1$

$$\frac{\partial J}{\partial z_i} = \hat{y}_i - y_i$$

General Delta Rule

$$\frac{\partial J}{\partial w_{ij}} = \delta_i a_j$$



# General Supervised Learning Algorithm

- Step 1 : Data Collection – Collect Data in  $(x, y)$  pairs
- Step 2 : Parameter Initialization – Guess for  $w$
- Step 3 : Forward Pass – For all data points in the dataset find  $\hat{y}$  using a hypothesis Model
- Step 4 : Gradient Descent – Improve  $w$  using gradient descent.  $w_j = w_j - \alpha \frac{\partial J}{\partial w_j}$
- Step 5 : Check exit criterion.
  - If not satisfied, go to Step 3 with the updated parameters

$$\mathbf{x}^{(1)} = (x_1^{(1)}, x_2^{(1)}, \dots, x_n^{(1)})$$

$$\mathbf{y}^{(1)} = (y_1^{(1)}, y_2^{(1)}, \dots, y_K^{(1)})$$

$$\hat{\mathbf{y}}^{(1)} = (\hat{y}_1^{(1)}, \hat{y}_2^{(1)}, \dots, \hat{y}_K^{(1)})$$

x	y	$\hat{y}$	J
$\mathbf{x}^{(1)}$	$\mathbf{y}^{(1)}$	$\hat{\mathbf{y}}^{(1)}$	$J^{(1)}$
$\mathbf{x}^{(2)}$	$\mathbf{y}^{(2)}$	$\hat{\mathbf{y}}^{(2)}$	$J^{(2)}$
$\mathbf{x}^{(3)}$	$\mathbf{y}^{(3)}$	$\hat{\mathbf{y}}^{(3)}$	$J^{(3)}$
$\mathbf{x}^{(4)}$	$\mathbf{y}^{(4)}$	$\hat{\mathbf{y}}^{(4)}$	$J^{(4)}$
$\mathbf{x}^{(5)}$	$\mathbf{y}^{(5)}$	$\hat{\mathbf{y}}^{(5)}$	$J^{(5)}$
...	...	...	...
$\mathbf{x}^{(m)}$	$\mathbf{y}^{(m)}$	$\hat{\mathbf{y}}^{(m)}$	$J^{(m)}$
			$J = \sum J^{(i)}$

# Linear Regression – Learning Algorithm

- Step 1 : Data Collection – Collect Data in  $(x, y)$  pairs
- Step 2 : Parameter Initialization – Guess for  $w$
- Step 3 : Forward Pass – For all data points in the dataset  
find  $\hat{y}$  using
  - The Linear Model  $y^{(j)} = \sum_{i=0}^n w_i x_i^{(j)}$
- Step 4 : Gradient Descent – Improve  $w$  using gradient descent.  $w_j = w_j - \alpha \frac{\partial J}{\partial w_j}$ 
  - Gradient  $\frac{\partial J}{\partial w_j} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$
- Step 5 : Check exit criterion.
  - If not satisfied, go to Step 3 with the updated parameters

$x$	$y$	$\hat{y}$	$J$
$x^{(1)}$	$y^{(1)}$	$\hat{y}^{(1)}$	$J^{(1)}$
$x^{(2)}$	$y^{(2)}$	$\hat{y}^{(2)}$	$J^{(2)}$
$x^{(3)}$	$y^{(3)}$	$\hat{y}^{(3)}$	$J^{(3)}$
$x^{(4)}$	$y^{(4)}$	$\hat{y}^{(4)}$	$J^{(4)}$
$x^{(5)}$	$y^{(5)}$	$\hat{y}^{(5)}$	$J^{(5)}$
...	...	...	...
$x^{(m)}$	$y^{(m)}$	$\hat{y}^{(m)}$	$J^{(m)}$
			$J = \sum J^{(i)}$

# Binary Classification – Learning Algorithm

- Step 1 : Data Collection – Collect Data in  $(x, y)$  pairs
- Step 2 : Parameter Initialization – Guess for  $w$
- Step 3 : Forward Pass – For all data points in the dataset find  $\hat{y}$  using
  - The Logistic Model  $\hat{y}^{(j)} = \sigma(\sum_{i=0}^n w_i x_i^{(j)})$
- Step 4 : Gradient Descent – Improve  $w$  using gradient descent.  $w_j = w_j - \alpha \frac{\partial J}{\partial w_j}$ 
  - Gradient  $\frac{\partial J}{\partial w_j} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$
- Step 5 : Check exit criterion.
  - If not satisfied, go to Step 3 with the updated parameters

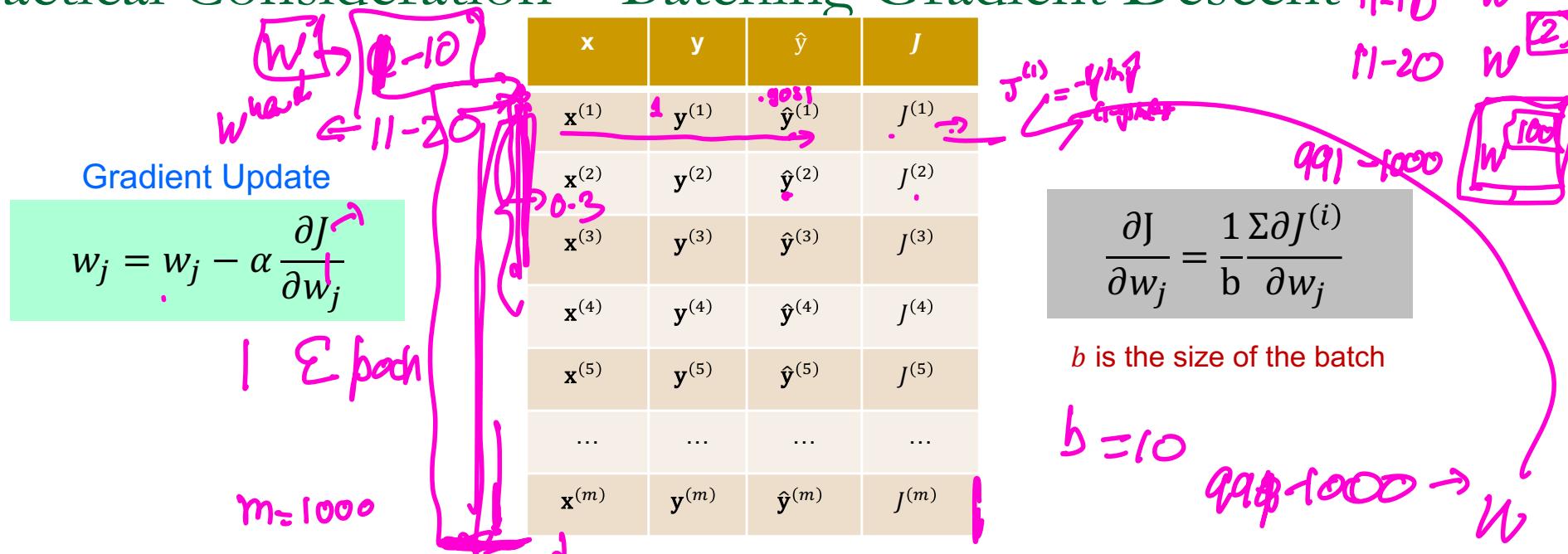
$x$	$y$	$\hat{y}$	$J$
$x^{(1)}$	$y^{(1)}$	$\hat{y}^{(1)}$	$J^{(1)}$
$x^{(2)}$	$y^{(2)}$	$\hat{y}^{(2)}$	$J^{(2)}$
$x^{(3)}$	$y^{(3)}$	$\hat{y}^{(3)}$	$J^{(3)}$
$x^{(4)}$	$y^{(4)}$	$\hat{y}^{(4)}$	$J^{(4)}$
$x^{(5)}$	$y^{(5)}$	$\hat{y}^{(5)}$	$J^{(5)}$
...	...	...	...
$x^{(m)}$	$y^{(m)}$	$\hat{y}^{(m)}$	$J^{(m)}$
			$J = \sum J^{(i)}$

# Multinomial Classification – Learning Algorithm

- Step 1 : Data Collection – Collect Data in  $(x, y)$  pairs
- Step 2 : Parameter Initialization – Guess for  $w$
- Step 3 : Forward Pass – For all data points in the dataset find  $\hat{y}$  using
  - The Softmax Model  $\hat{y}^{(j)} = \text{softmax}(\sum_{i=0}^n w_i x_i^{(j)})$
- Step 4 : Gradient Descent – Improve  $w$  using gradient descent.  $w_j = w_j - \alpha \frac{\partial J}{\partial w_j}$ 
  - Gradient  $\frac{\partial J}{\partial w_j} = \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$
- Step 5 : Check exit criterion.
  - If not satisfied, go to Step 3 with the updated parameters

$x$	$y$	$\hat{y}$	$J$
$x^{(1)}$	$y^{(1)}$	$\hat{y}^{(1)}$	$J^{(1)}$
$x^{(2)}$	$y^{(2)}$	$\hat{y}^{(2)}$	$J^{(2)}$
$x^{(3)}$	$y^{(3)}$	$\hat{y}^{(3)}$	$J^{(3)}$
$x^{(4)}$	$y^{(4)}$	$\hat{y}^{(4)}$	$J^{(4)}$
$x^{(5)}$	$y^{(5)}$	$\hat{y}^{(5)}$	$J^{(5)}$
...	...	...	...
$x^{(m)}$	$y^{(m)}$	$\hat{y}^{(m)}$	$J^{(m)}$
			$J = \sum J^{(i)}$

# Practical Consideration – Batching Gradient Descent



- Default option – Each update in gradient descent requires to see the whole dataset
  - That is, one update requires us to average the updates of every point in the whole dataset
  - This is called **Batch Gradient Descent** where  $b = m$ . **One epoch has one update**
  - Problems : (i) It is slow (ii) It is expensive for each update (iii) Data cannot be dynamic
- **Stochastic Gradient Descent**
  - $b = 1$ . **One epoch has  $m$  updates**. Many more updates per epoch
  - Can handle dynamic incoming data but loss convergence is more erratic
- **Mini-Batch Gradient Descent**
  - $1 < b < m$ . **One epoch has  $\frac{m}{b}$  updates**. More than 1 update per epoch
  - Loss convergence more erratic than Batch and less erratic than stochastic

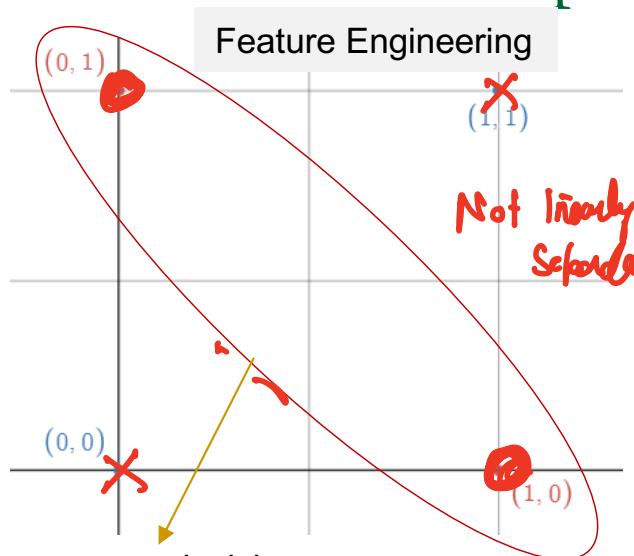
# The XOR problem – Need for deep networks

	$x_1$	$x_2$	$y$
1	0	0	0
2	1	0	1
3	0	1	1
4	1	1	0

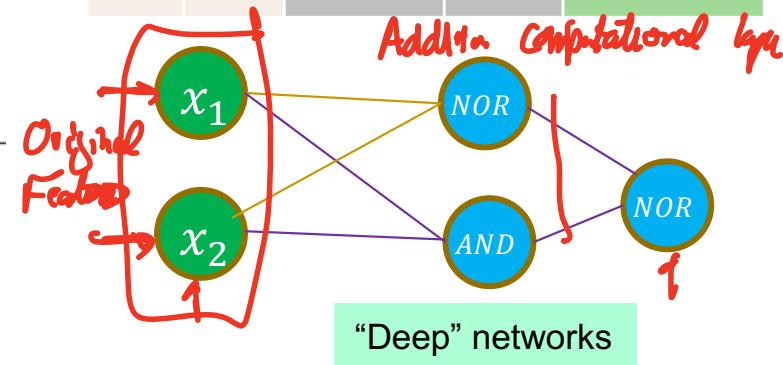
$$\hat{y} = x_1 + x_2 - 2x_1x_2$$

hand-picked

Additional factors  $\rightarrow$  Feature learning!



$x_1$	$x_2$	NOR	AND	$NOR(NOR, AND)$
0	0	1	0	0
1	0	0	0	1
0	1	0	0	1
1	1	0	1	0



- In all cases so far, we map the input to  $z = w_0 + w_1x_1 + \dots$
- The decision boundary,  $z = 0$  will, therefore, always be a line
  - If we directly map  $x \rightarrow z$  like so far, it only works for **linearly separable data**
- Two solutions
  - Feature engineering -- Use nonlinear features. Example --  $\hat{y} = w_0 + w_1x_1 + w_2x_2 + w_3x_1x_2$
  - Use “deep” networks – Add extra hidden layers to create non-interpretable nonlinear features
- We can use  $XOR(x_1, x_2) = NOR(NOR(x_1, x_2), AND(x_1, x_2))$  for example
- But, can we automatically come up with this representation via learning?

$$x_3 = x_1x_2 \text{ Nonlinear feature}$$

Physics Informed  
Neural Networks

$$V = u + gt$$

PINN

Raissi & Karniadakis



## DEEP NEURAL NETWORKS

(Multi Layer Perceptrons – MLPs)

$$\text{Regression tasks} \rightarrow g_0(z) = z \quad Z = w_1a_1 + w_2a_2 + v_3a_3 + v_4a_4 + v_5a_5 + w_0$$

Linear activation  $\hat{y} = g_0(z)$  Nonlinearity at the output

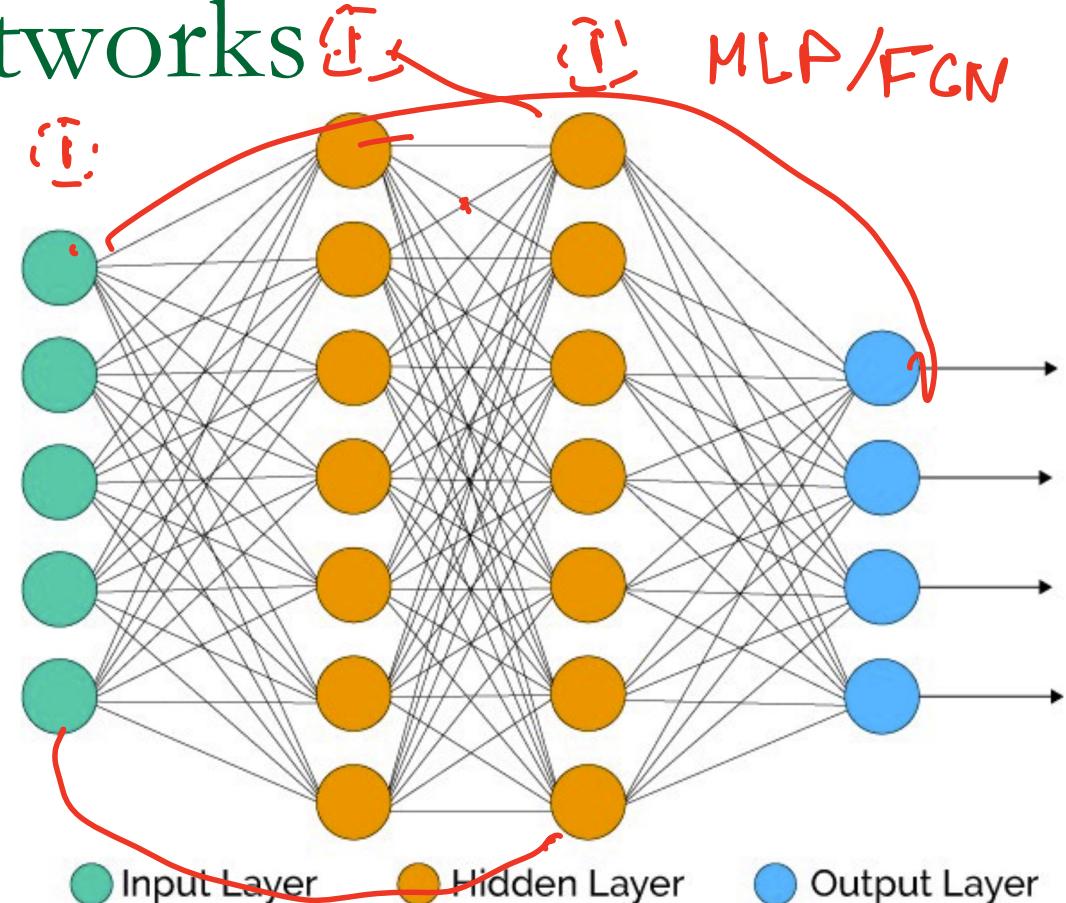
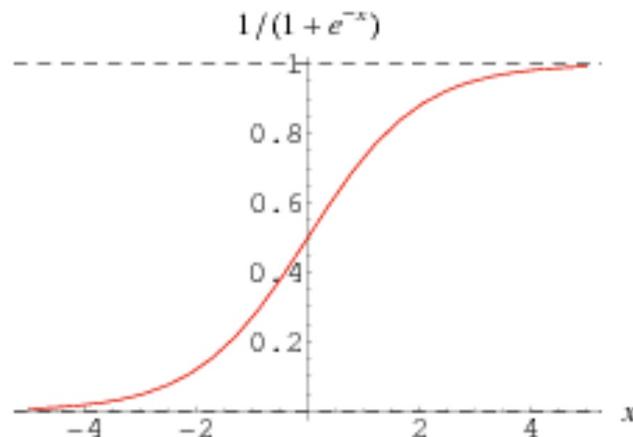
Can be different from  $g$  inside the network

# Recall : Neural Networks

## Neural Network Model

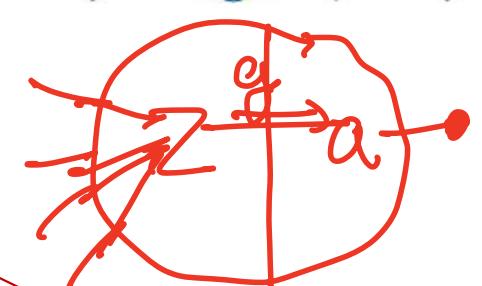
Common Activation function

$$\text{Sigmoid}(x) = \frac{1}{1+\exp(-x)}$$

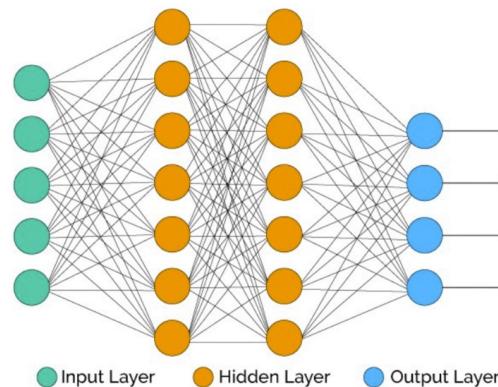
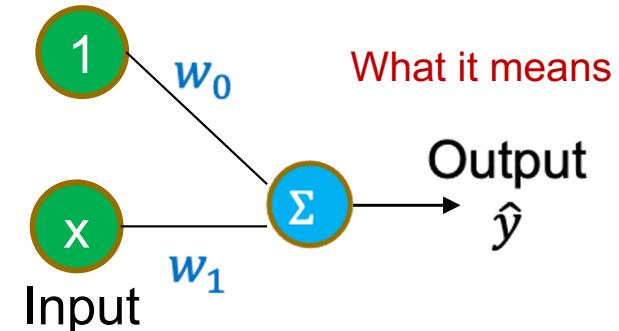
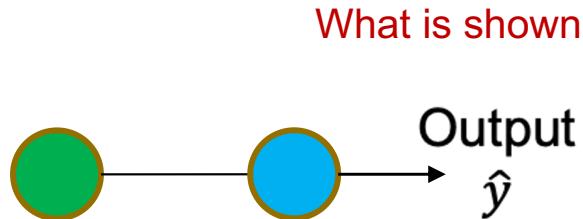


$$a = g(\sum w_i x_i)$$

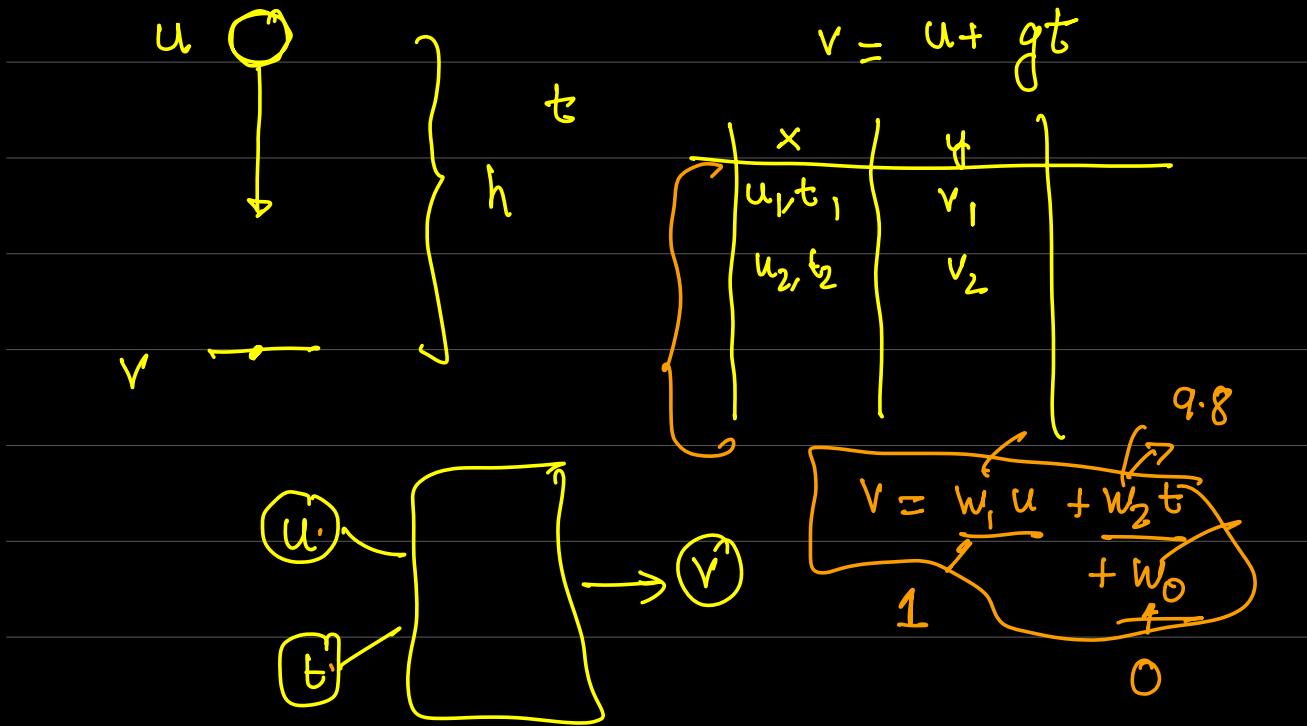
A Neuron has a linear and a nonlinear operation



# Weights and biases

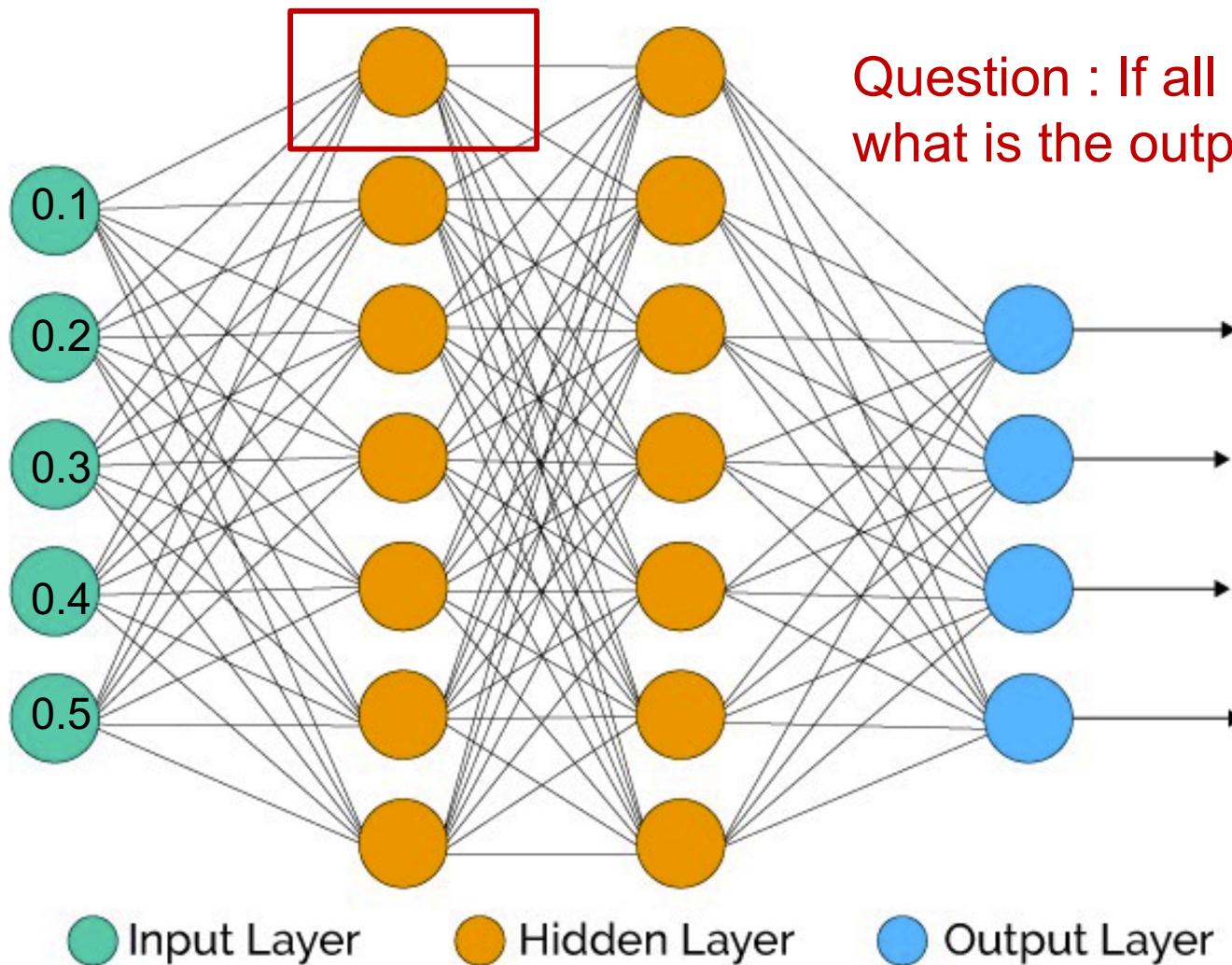


- Constant neurons representing a “feature” with 1 are usually not depicted in the pictorial representation.
- The parameters going from the constant neuron is called a bias and the neuron is called a “bias unit”
- The other parameters are called “weights”

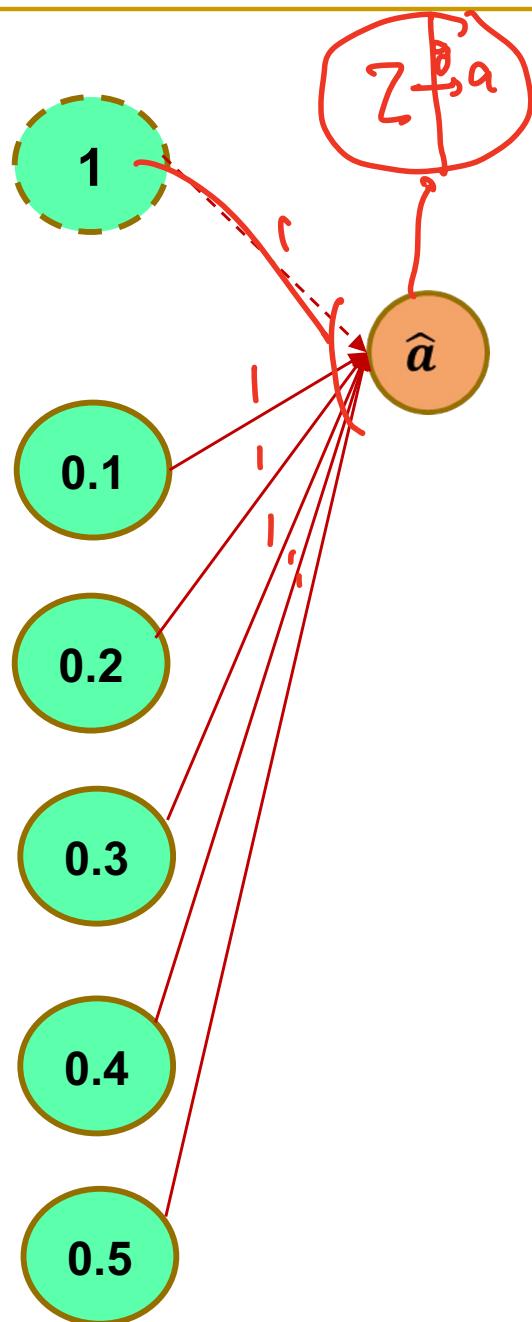


$$h = \frac{v^2 - u^2}{2g}$$

# The output of a single neuron



Question : If all weights are 1, then what is the output of the shown neuron?



Find the output of the neuron  $\hat{a}$

All weights are equal to 1

**Ans:** The output has two parts

**Linear** --  $z = \sum w_i x_i$

$$z = w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5$$

$$z = 1 + 0.1 + 0.2 + 0.3 + 0.4 + 0.5 = 2.5$$

**Nonlinear Activation** --  $\hat{a} = g(z)$

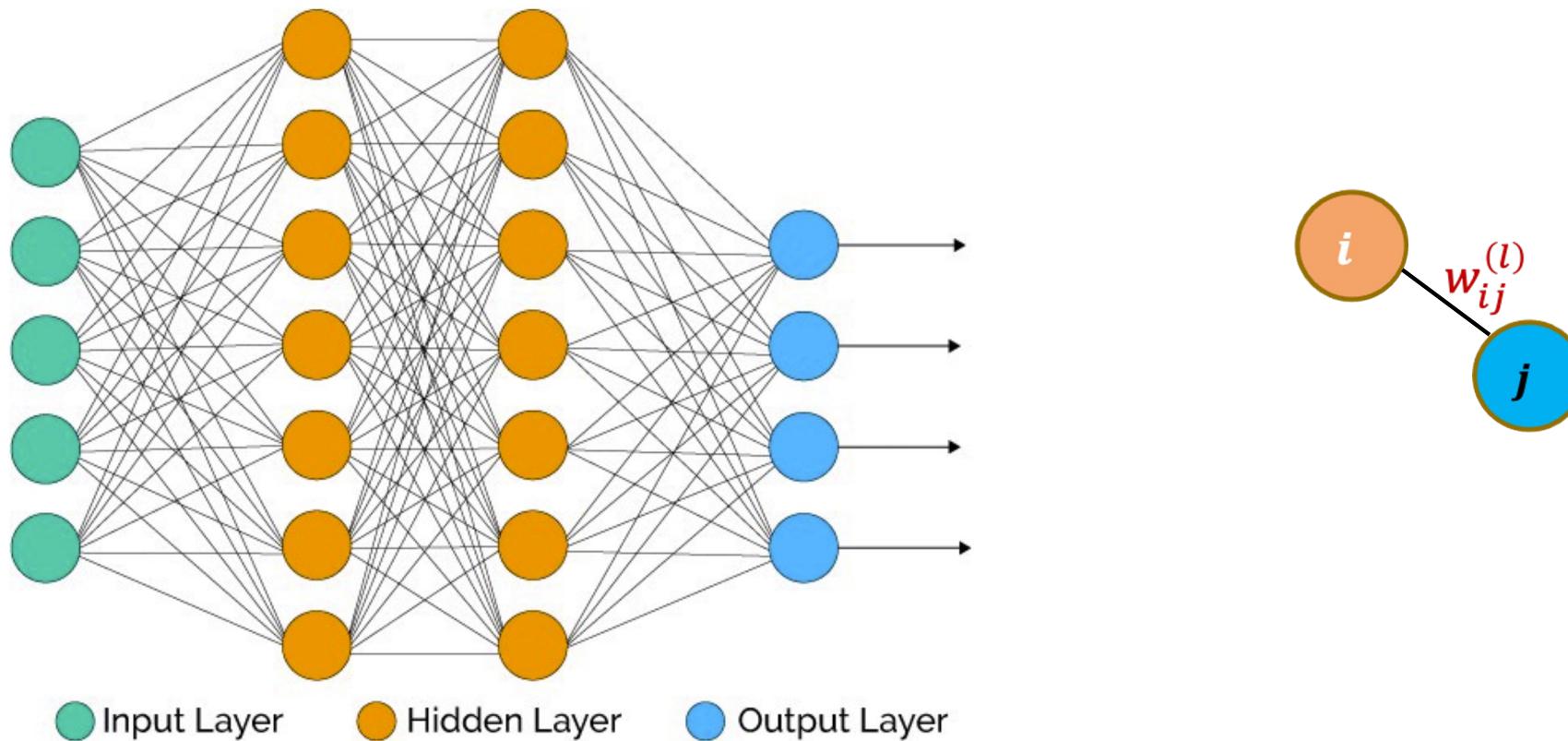
$$= \frac{1}{(1 + \exp(-2.5))}$$

$$= 0.9241$$

**Don't forget the bias unit!**

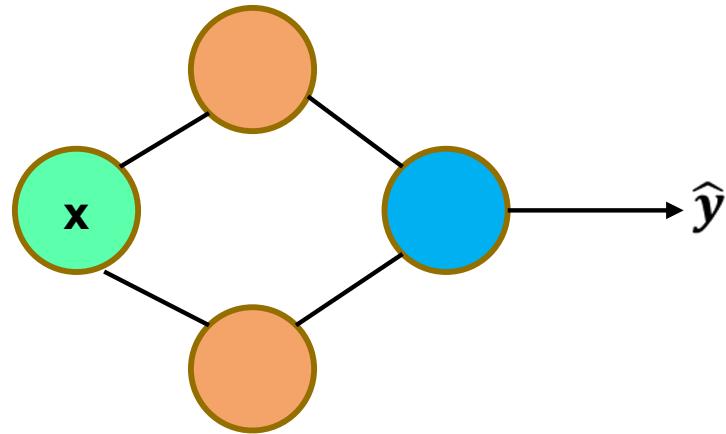
# Weights notation

(temporarily, Only for this example)

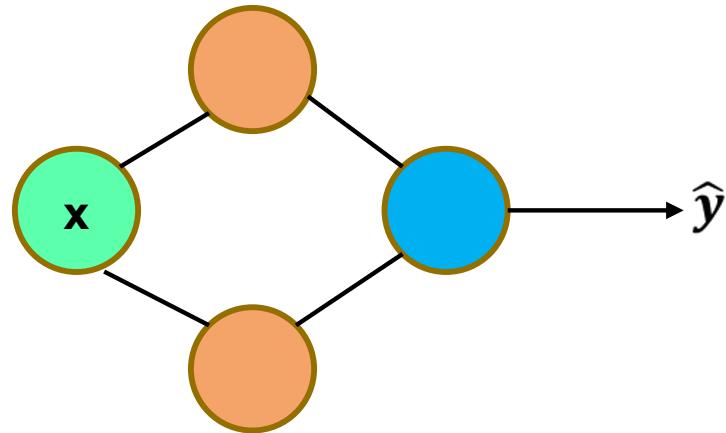


- $w_{ij}^{(l)}$  is the weight connecting Neuron  $i$  of Layer  $l$  to Neuron  $j$  of Layer  $l + 1$
- Bias units are the 0<sup>th</sup> neurons

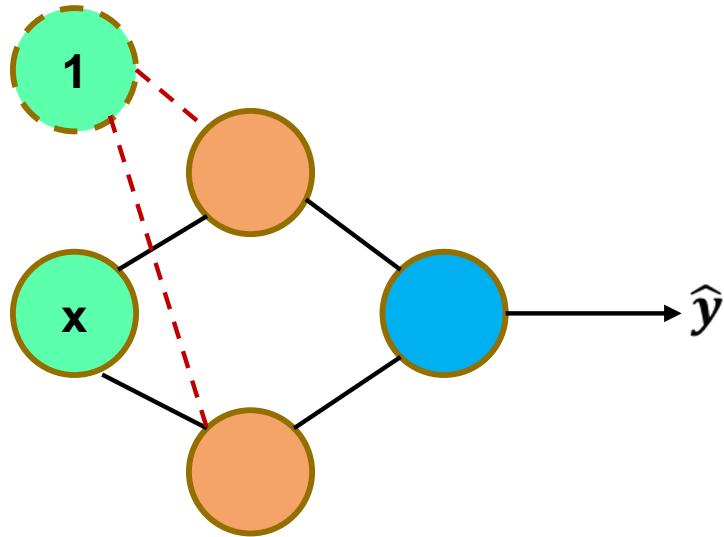
A simple forward pass calculation



**Question:** Find the output of the given network for the following



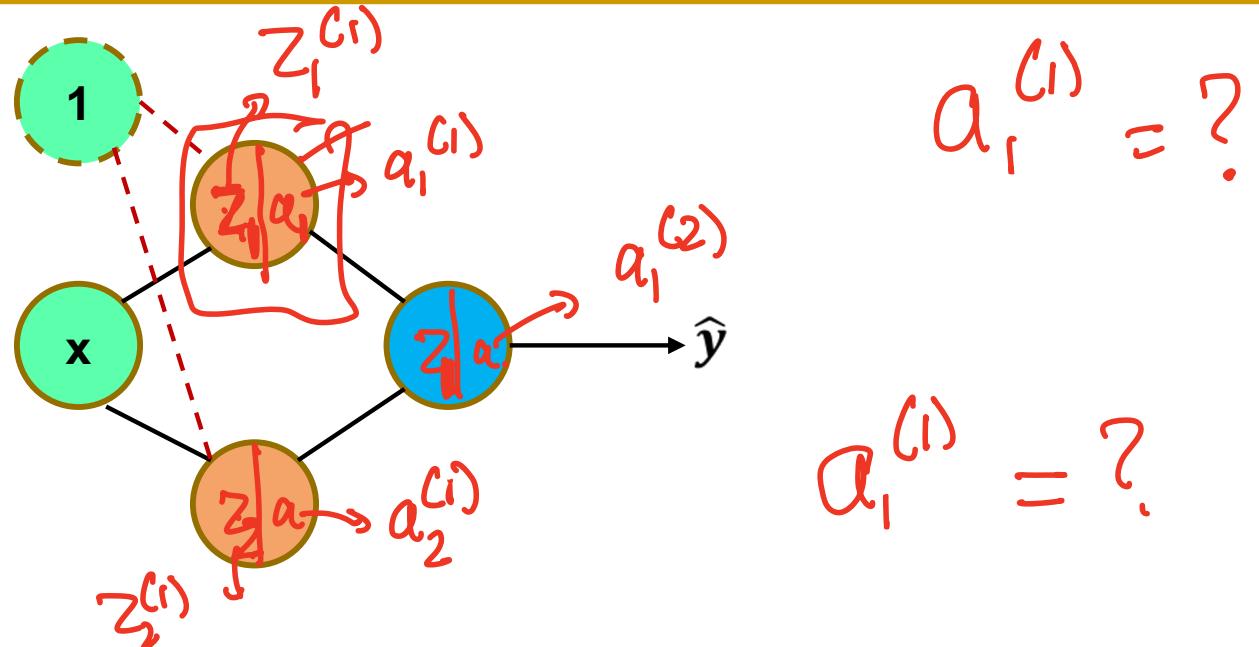
**Question:** Find the output of the given network for the following  
 $x = 0.5$



**Question:** Find the output of the given network for the following

$$x = 0.5$$

$$w_{01}^{(1)} = 1.0 \quad w_{02}^{(1)} = 0.8$$



$$a_1^{(1)} = ?$$

$$a_1^{(1)} = ?$$

**Question:** Find the output of the given network for the following

$$x = 0.5$$

$$w_{01}^{(1)} = 1.0 \quad w_{02}^{(1)} = 0.8$$

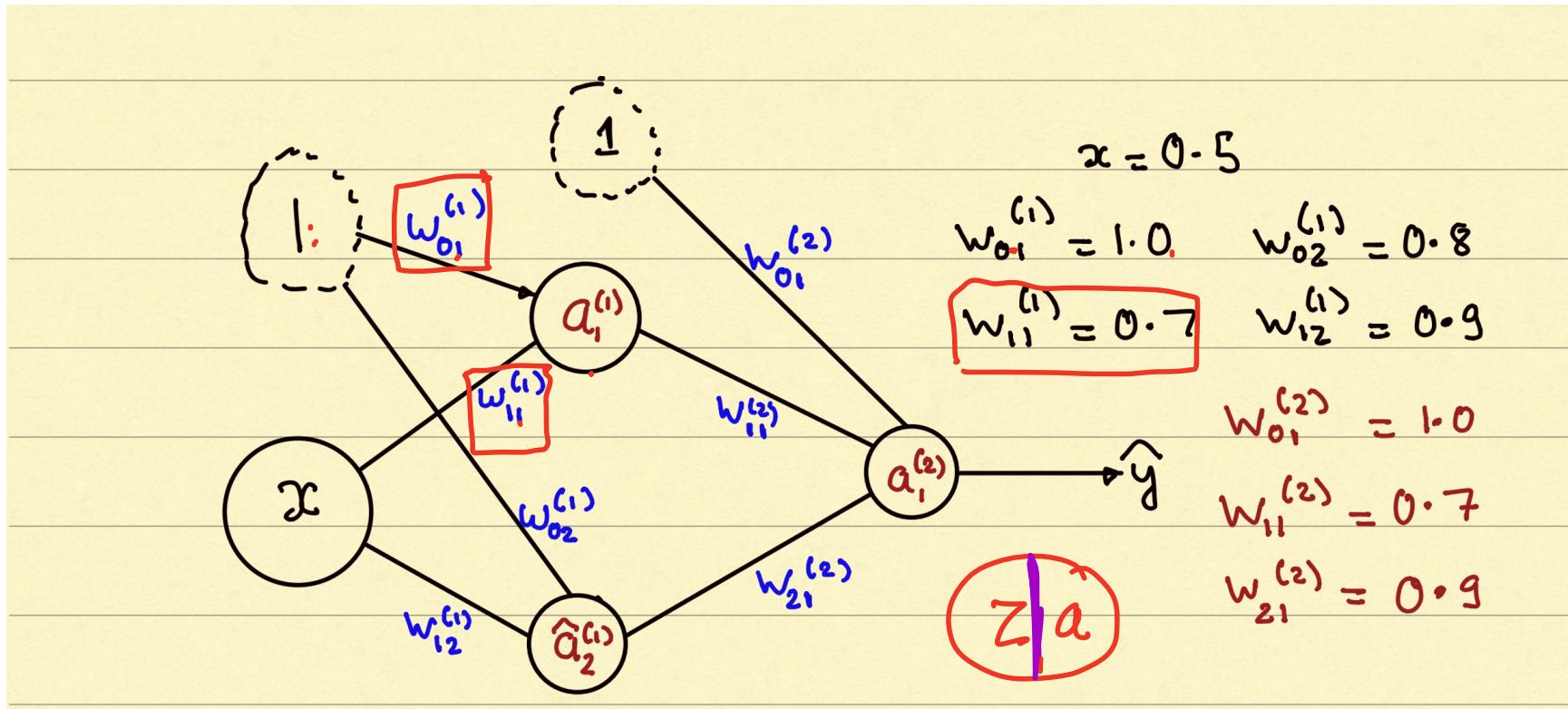
$$w_{11}^{(1)} = 0.7 \quad w_{12}^{(1)} = 0.9$$

$$w_{01}^{(2)} = 1$$

$$w_{11}^{(2)} = 0.7 \quad w_{21}^{(2)} = 0.9$$

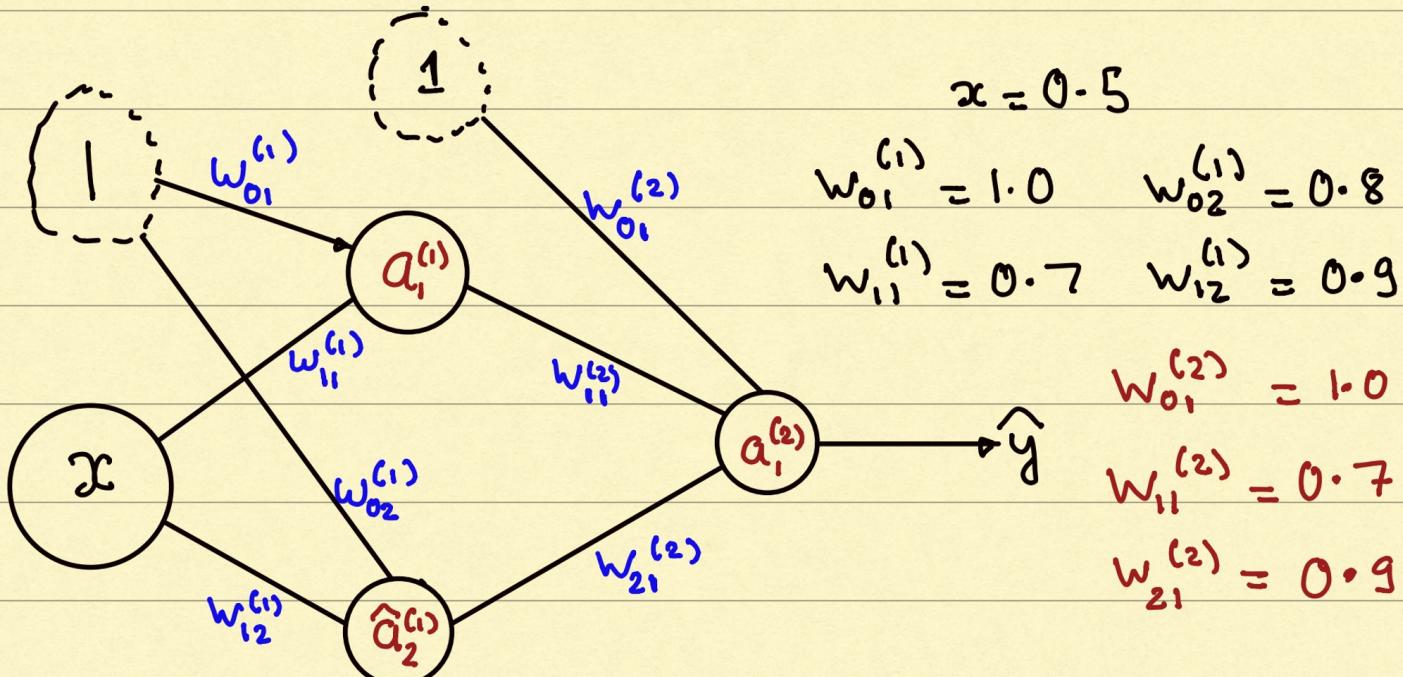
Assume that the activation function in the hidden layer and the output layer is the sigmoid

# SOLUTION

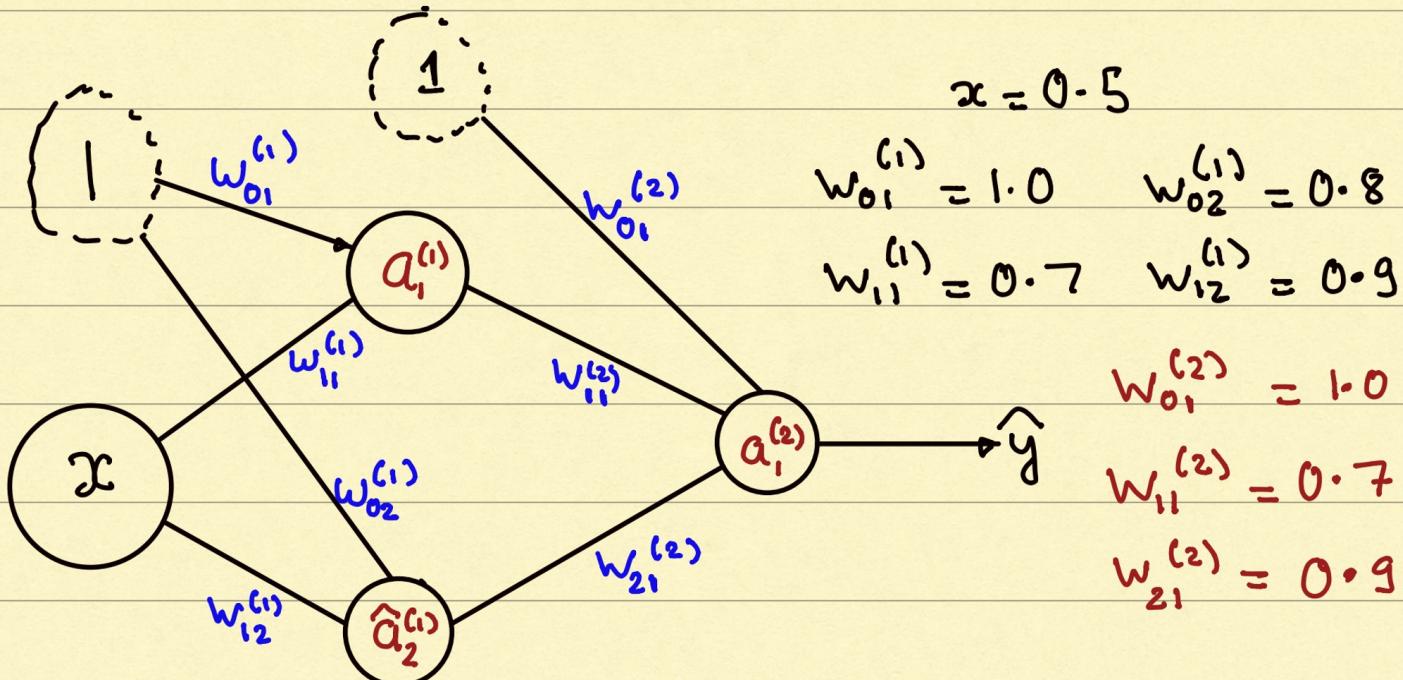


$$\begin{aligned}
 \text{Linear} \rightarrow Z_1^{(1)} &= w_{01}^{(1)} * 1 + w_{11}^{(1)} * x \\
 &= 1 * 1 + 0 \cdot 7 * 0 \cdot 5 \\
 &= 1 \cdot 35
 \end{aligned}$$

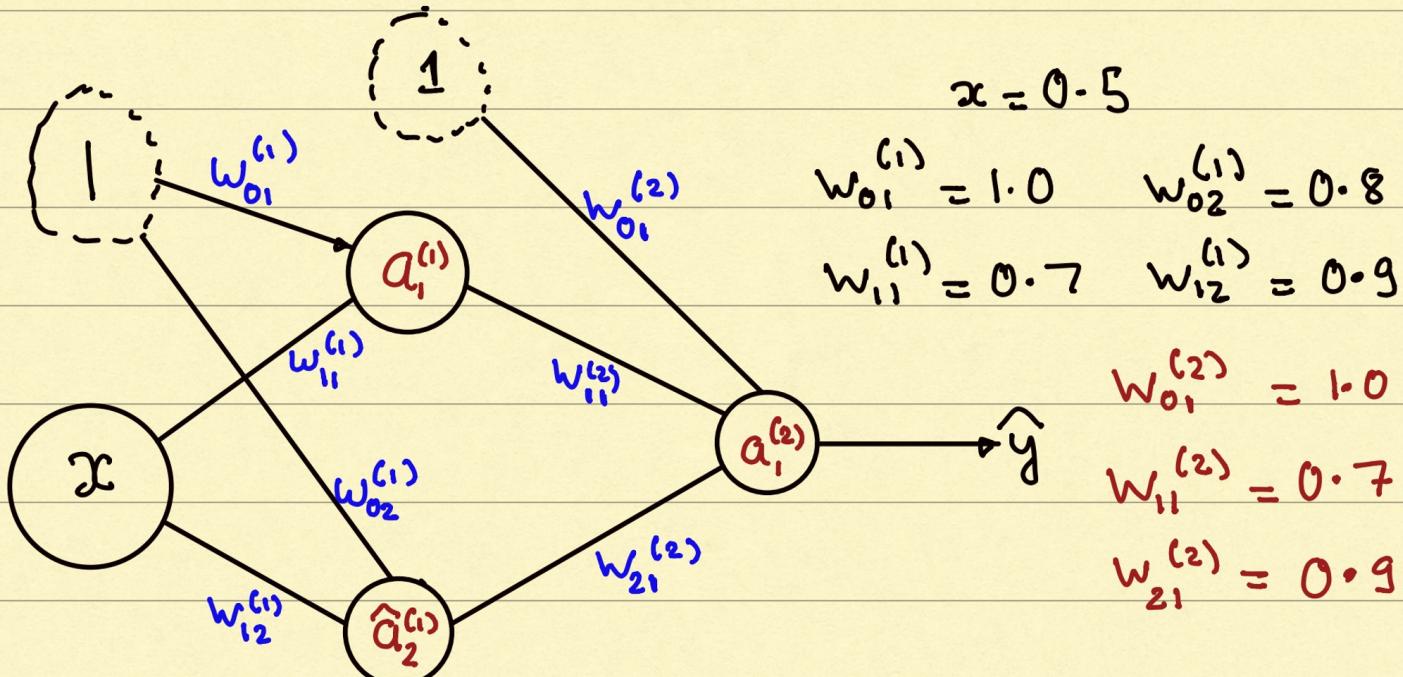
$$\begin{aligned}
 a_1^{(1)} &= g(Z_1^{(1)}) \\
 &= \sigma(1 \cdot 35)
 \end{aligned}$$



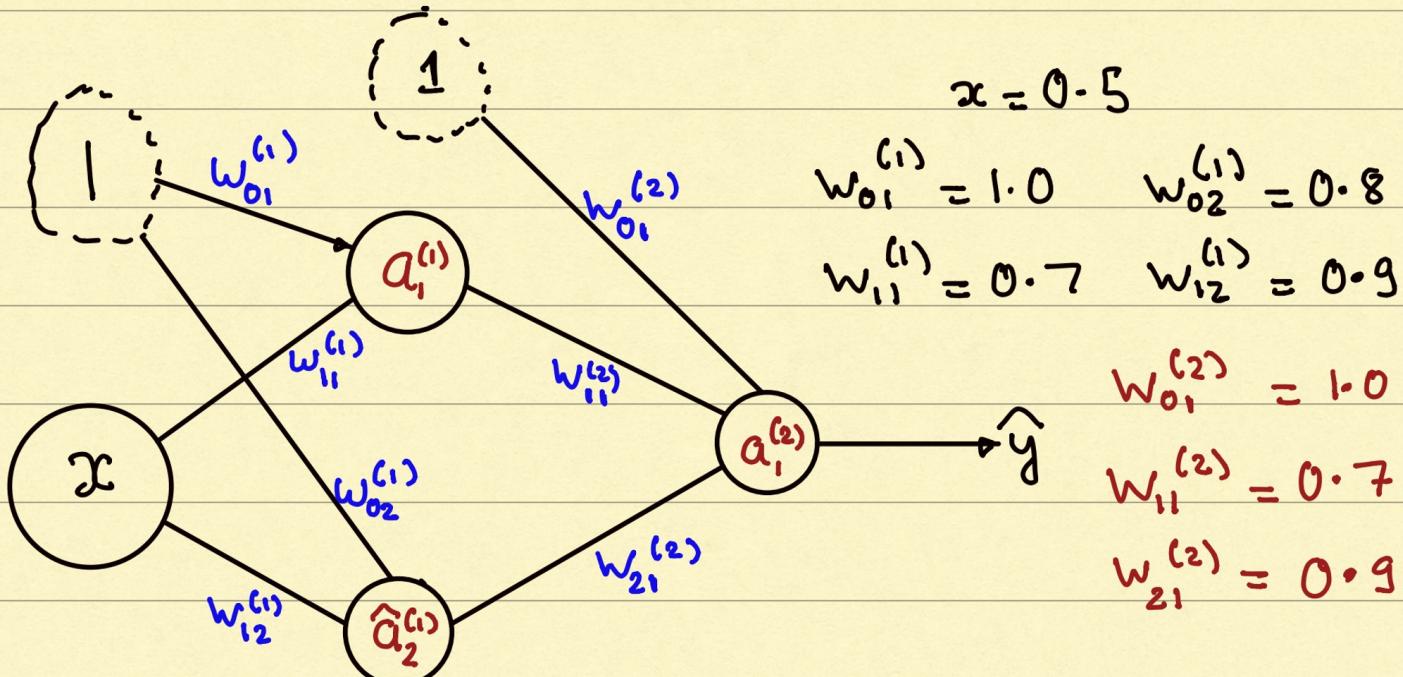
$$\hat{a}_1^{(1)} = g(z_1^{(1)}) =$$



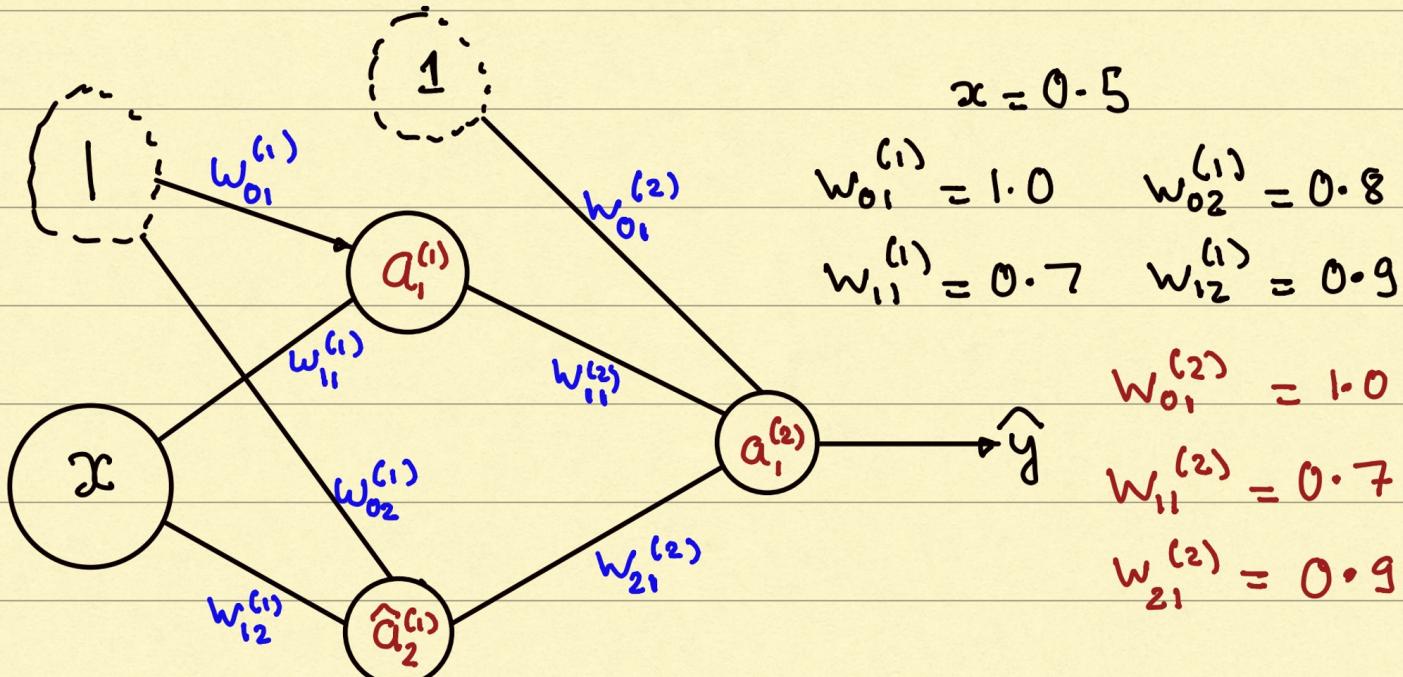
$$\hat{a}_1^{(1)} = g(z_1^{(1)}) = g(w_{01}^{(1)}x_0 + w_{11}^{(1)}x)$$



$$\begin{aligned}
 \hat{a}_1^{(1)} &= g(z_1^{(1)}) = g(w_{01}^{(1)}x_0 + w_{11}^{(1)}x) \\
 &= g(1*1 + 0.7*0.5) = g(1.35)
 \end{aligned}$$

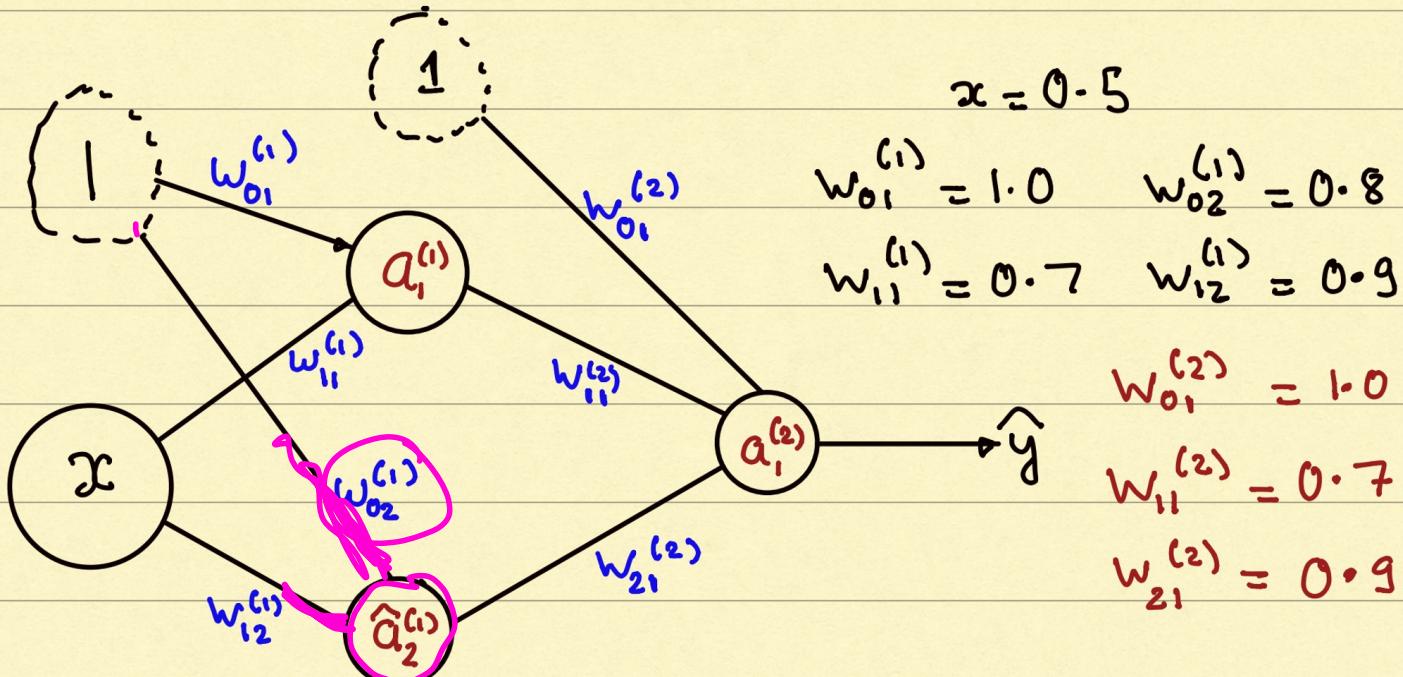


$$\begin{aligned}
 \hat{a}_1^{(1)} &= g(z_1^{(1)}) = g(w_{01}^{(1)}x_0 + w_{11}^{(1)}x) \\
 &= g(1 \cdot 1 + 0 \cdot 7 \cdot 0 \cdot 5) = g(1 \cdot 35)
 \end{aligned}$$



$$\begin{aligned}
 \hat{a}_1^{(1)} &= g(z_1^{(1)}) = g(w_{01}^{(1)}x_0 + w_{11}^{(1)}x) \\
 &= g(1*1 + 0.7*0.5) = g(1.35)
 \end{aligned}$$

$\Rightarrow \hat{a}_1^{(1)} = 0.7941$



$$x = 0 \cdot 5$$

$$w_{01}^{(1)} = 1 \cdot 0 \quad w_{02}^{(1)} = 0 \cdot 8$$

$$w_{11}^{(1)} = 0 \cdot 7 \quad w_{12}^{(1)} = 0 \cdot 9$$

$$w_{01}^{(2)} = 1 \cdot 0$$

$$w_{11}^{(2)} = 0 \cdot 7$$

$$w_{21}^{(2)} = 0 \cdot 9$$

$$a_2^{(1)} \rightarrow$$

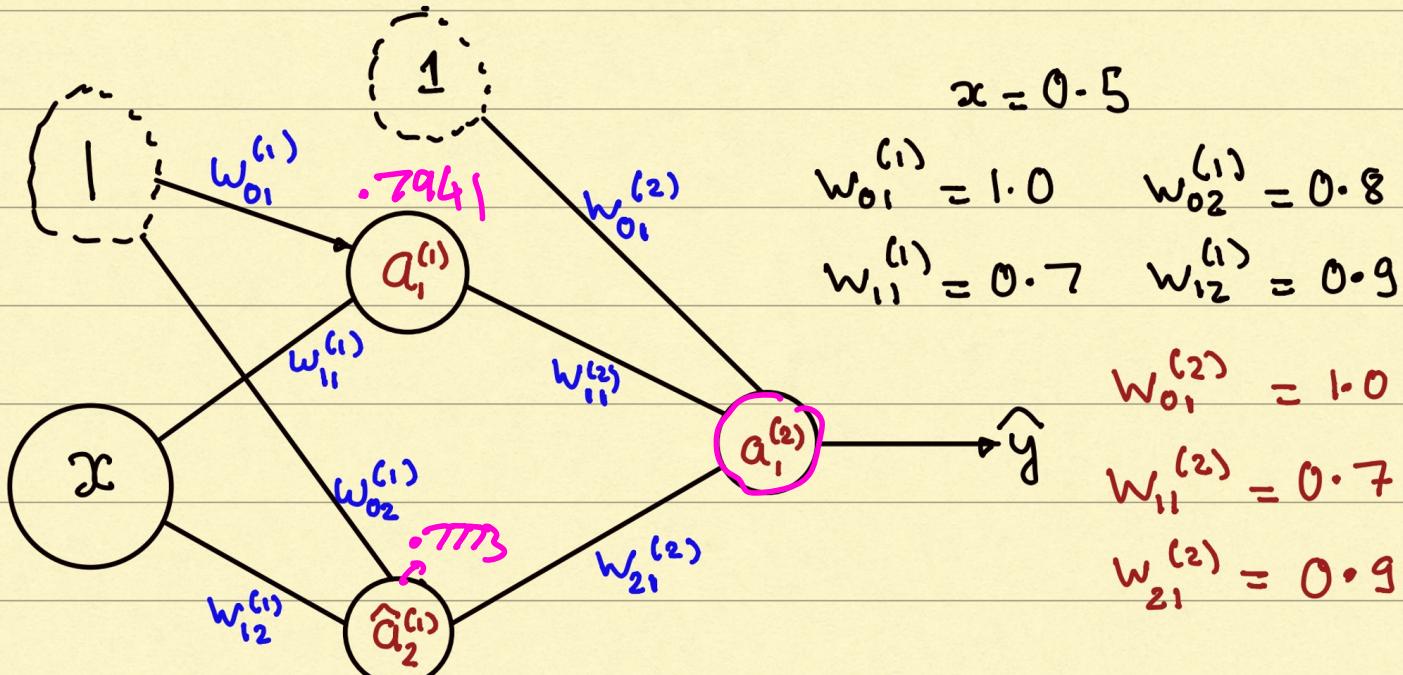


$$z_2^{(1)} = w_{02}^{(1)} * 1 + w_{12}^{(1)} * 0.5$$

$$0.8 * 1 + 0.9 * 0.5$$

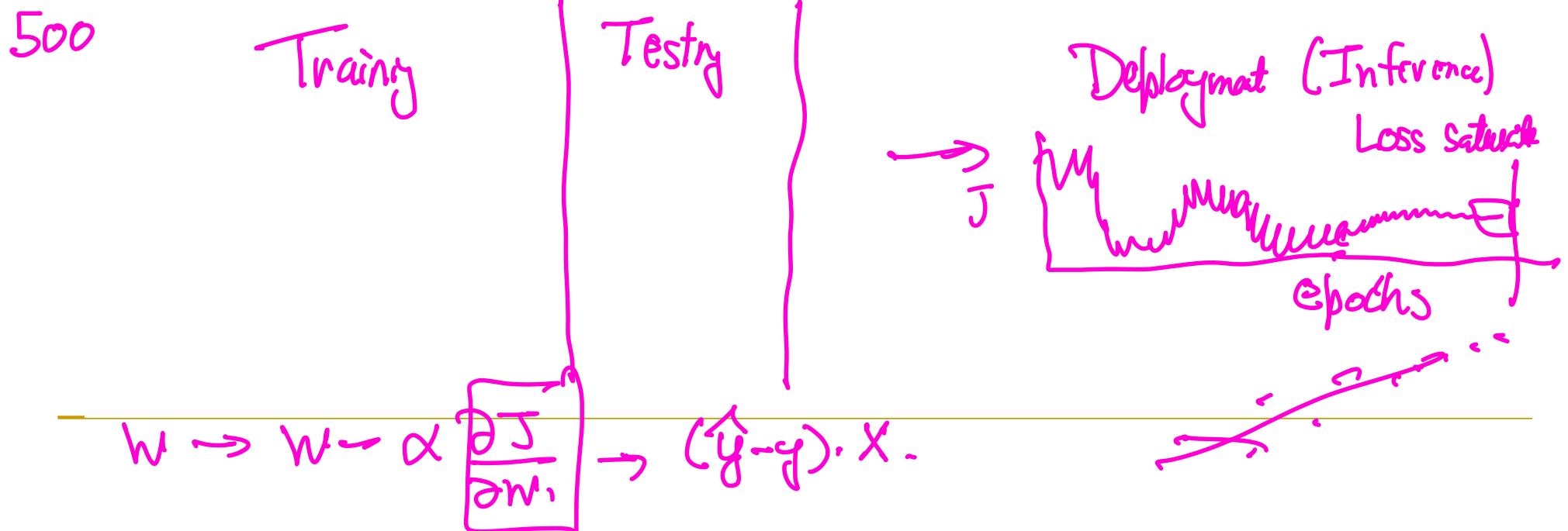
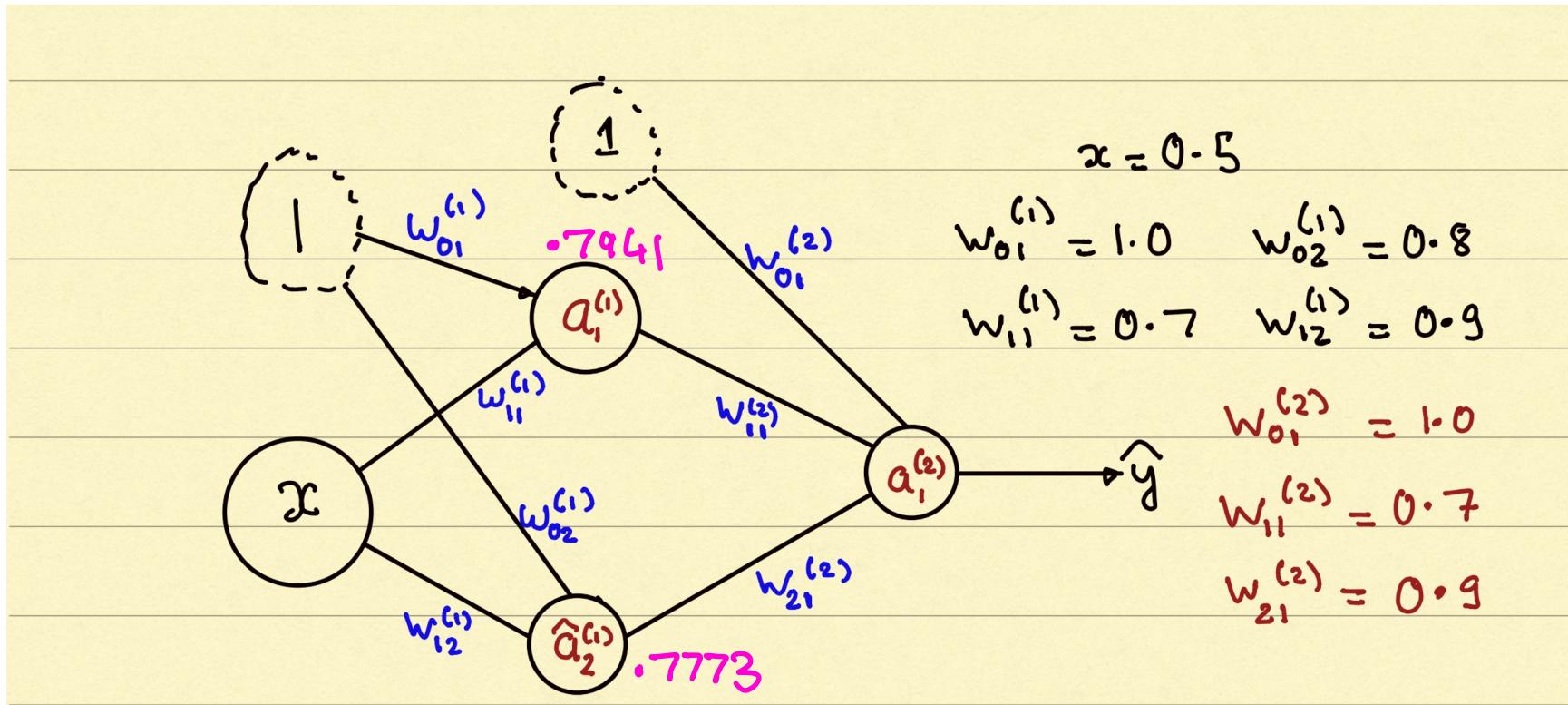
$$= 0.8 + 0.45 = 1.25$$

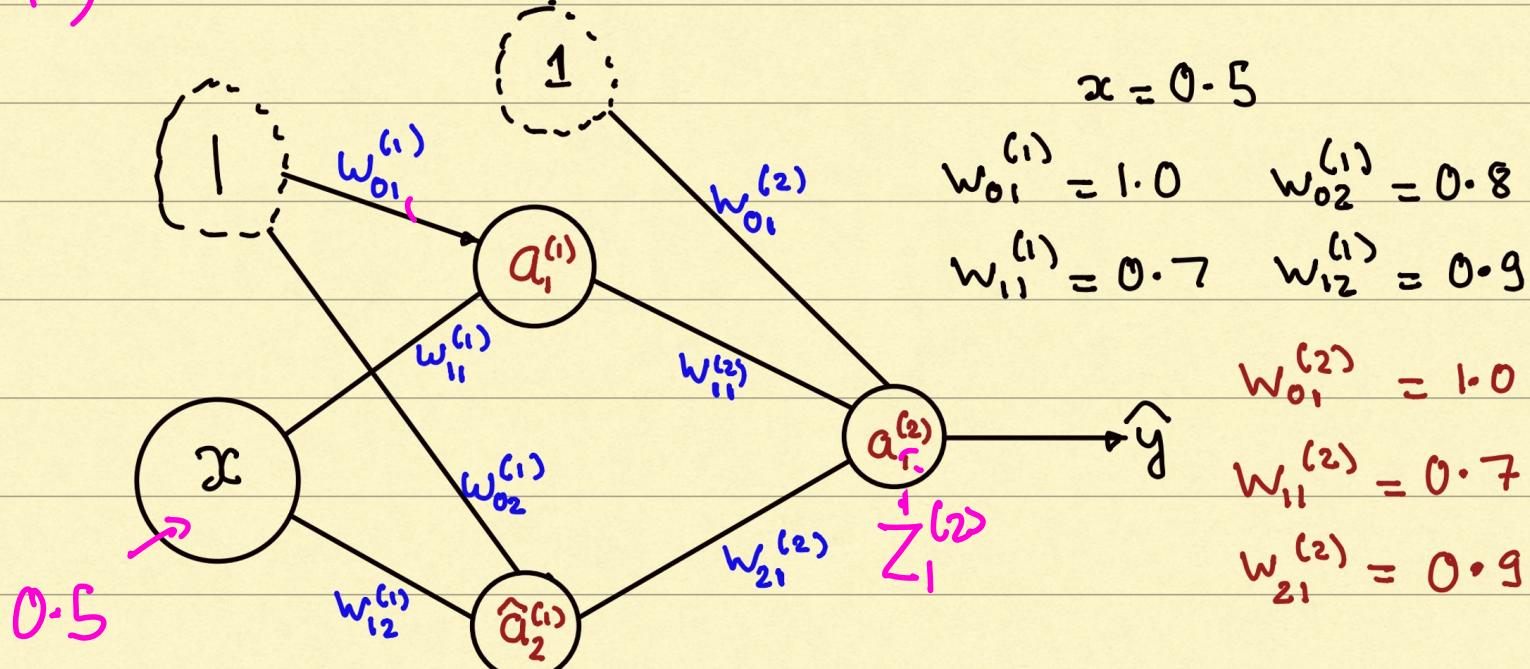
$$a_2^{(1)} = \sigma(1.25)$$



$$\begin{aligned}
 \hat{a}_2^{(1)} &= g(z_2^{(1)}); \quad z_2^{(1)} = w_{02}^{(1)} x_0 + w_{12}^{(1)} x_1 = 0.8 + 0.9 * 0.5 \\
 \Rightarrow z_2^{(1)} &= 1.25 \Rightarrow \hat{a}_2 = g(z_2^{(1)}) \\
 \Rightarrow \hat{a}_2^{(1)} &= 0.7773
 \end{aligned}$$

$$f(z) = \frac{1}{1 + \exp(-z)}$$





$$\hat{a}_1^{(2)} = g(\tilde{z}_1^{(2)}); \quad \tilde{z}_1^{(2)} = w_{01}^{(2)} x_0 + w_{11}^{(2)} a_1^{(1)} + w_{21}^{(2)} a_2^{(1)}$$

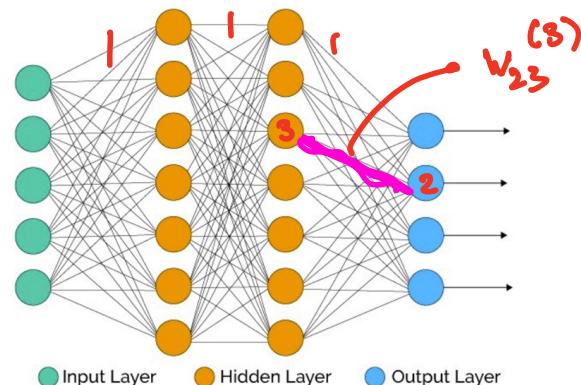
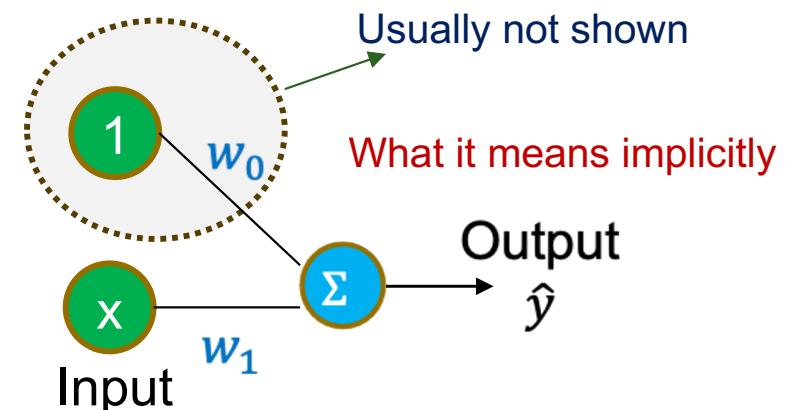
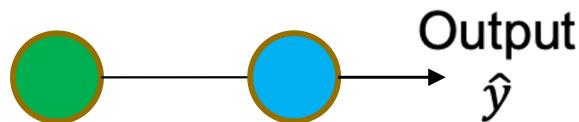
$$\Rightarrow \tilde{z}_1^{(2)} = 1 * 1 + 0.7 * 0.7941 + 0.9 * 0.7773$$

$$= 2.255 \Rightarrow a_1^{(2)} = 0.9051 = \hat{y}$$

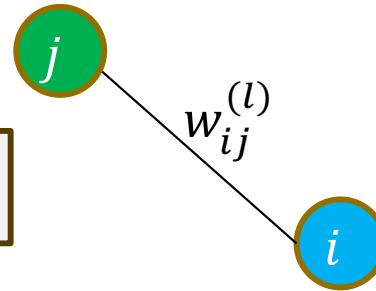
# Weights notation

(Actual  
Notation)

What is shown explicitly

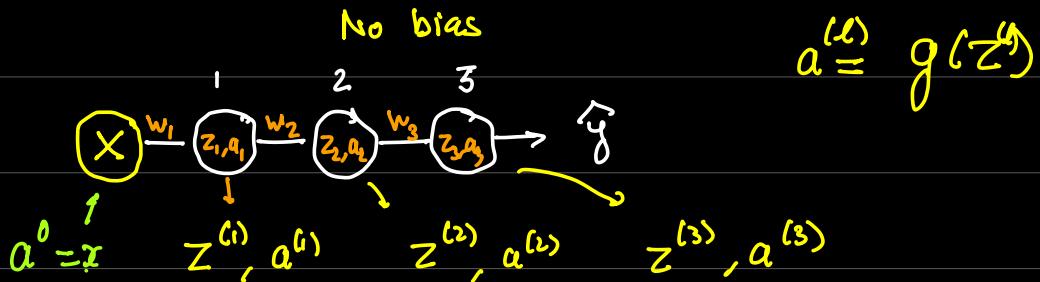


$w_{ij}^{(l)} \rightarrow w_{\text{to},\text{from}}^{(\text{layer})}$



- Constant neurons representing a “feature” with 1 are usually not depicted in the pictorial representation.
- The parameters going from the constant neuron is called a bias and the neuron is called a “bias unit”
- The other parameters are called “weights”

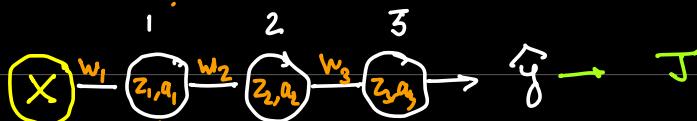
## Backprop in a scalar chain



Forward Pass —

$$\boxed{\begin{array}{ll} z^{(1)} = w_1 a^{(0)} & a^{(1)} = g(z^{(1)}) \\ z^{(2)} = w_2 a^{(1)} & a^{(2)} = g(z^{(2)}) \\ z^{(3)} = w_3 a^{(2)} & a^{(3)} = g(z^{(3)}) \end{array}}$$

$\hookrightarrow \hat{y} = a^{(3)}$



DELTA RULE

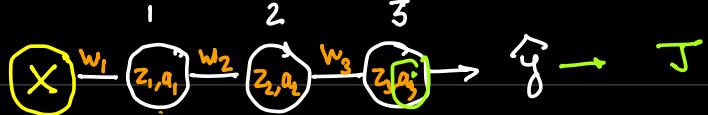
$\frac{\partial J}{\partial w_i}$	$\frac{\partial J}{\partial z_1} \equiv \delta_1 ; \quad \frac{\partial J}{\partial a_1}$
$\frac{\partial J}{\partial w_2}$	$\frac{\partial J}{\partial z_2} \equiv \delta_2 ; \quad \frac{\partial J}{\partial a_2}$
$\frac{\partial J}{\partial w_3}$	$\frac{\partial J}{\partial z_3} \equiv \delta_3 ; \quad \frac{\partial J}{\partial a_3}$



$$\begin{aligned} \frac{\partial J}{\partial w_i} &= \frac{\partial J}{\partial z^{(l+1)}} \cdot \frac{\partial z^{(l+1)}}{\partial w_i} \\ &= \delta^{(l+1)} a^{(l)} \end{aligned}$$

$$\underline{\underline{z}}^{(l+1)} = \underline{\underline{w}}_i a^{(l)}$$

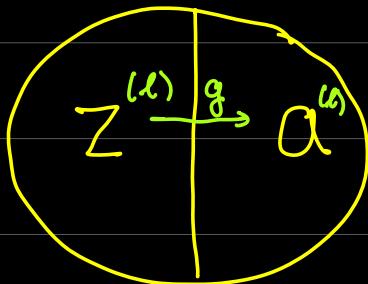
$$\boxed{\frac{\partial J}{\partial w_i} = \delta^{(l+1)} a^{(l)}}$$



$$\frac{\partial a}{\partial z} = g'(z)$$

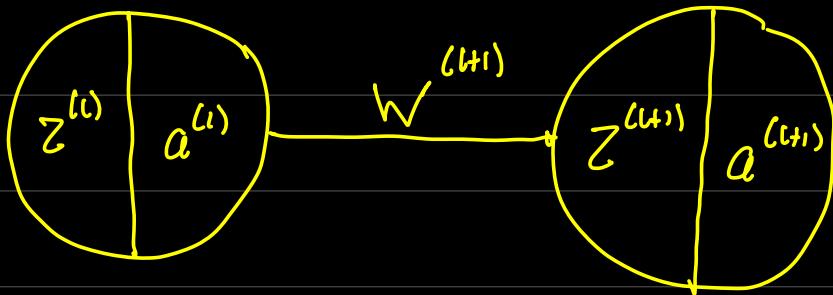
$$a = g(z)$$

$$\frac{\partial J}{\partial z^{(l)}}$$



$$\frac{\partial J}{\partial z^{(l)}} = \frac{\partial J}{\partial a^{(l)}} \frac{\partial a^{(l)}}{\partial z^{(l)}}$$

$$\delta^{(l)} = e^{(l)} \cdot \underbrace{g'(z^{(l)})}_{=}$$



$$\frac{\partial J}{\partial a^{(l)}} = \frac{\partial J}{\partial z^{(l+1)}} \cdot \frac{\partial z^{(l+1)}}{\partial a^{(l)}}$$

$$e^{(l)} = g^{(l+1)} \cdot w^{(l+1)}$$

$$\begin{cases} z^{(l+1)} = w^{(l+1)} a^{(l)} \\ g(z) = r(z) \\ g'(z) = \frac{e^z}{(1+e^{-z})} = \tau(z) \end{cases}$$

$$\frac{\partial J}{\partial w^{(l)}} = \delta^{(l)} a^{(l-1)}$$

$$e^{(l)} = g^{(l+1)} w^{(l+1)}$$

$$\begin{cases} \delta^{(l)} = e^{(l)} g'(z^{(l)}) \\ \tau(z) = r(z) \\ a^{(l-1)} = \underbrace{a^{(l)}}_{a^{(l-1)}} \end{cases}$$

Start with  $\frac{\partial J}{\partial a_3} = \frac{\partial J}{\partial \hat{y}} ; J = \frac{1}{2} (\hat{y} - y)^2$

$$\Rightarrow \frac{\partial J}{\partial \hat{y}} = \hat{y} - y$$

We know  $e_3$

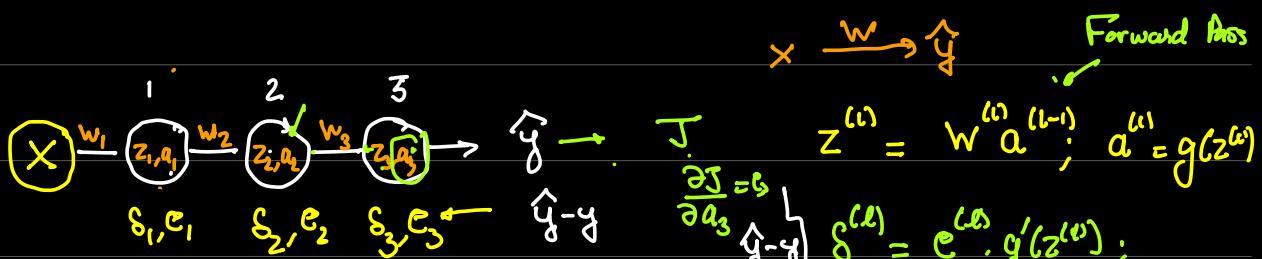
$$1-a = 1 - \frac{1}{1+e^{-z}}$$

$$\boxed{g(z) = \sigma(z)} = \frac{1}{1 + e^{-z}} = \frac{1 + e^{-z}}{e^{-z}} = \frac{1}{e^{-z}}$$

$$g'(z) = \frac{dg}{dz} = \frac{-1}{(1 + e^{-z})^2} * -e^{-z}$$

$$g'(z) = \frac{\frac{e^{-z}}{(1 + e^{-z})^2}}{1 + e^{-z}} = \frac{e^{-z}}{1 + e^{-z}} * \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-z}}$$

$$g'(z) = \boxed{a(1-a)} = \sigma(z)(1 - \sigma(z))$$

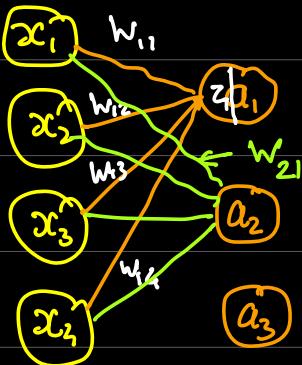
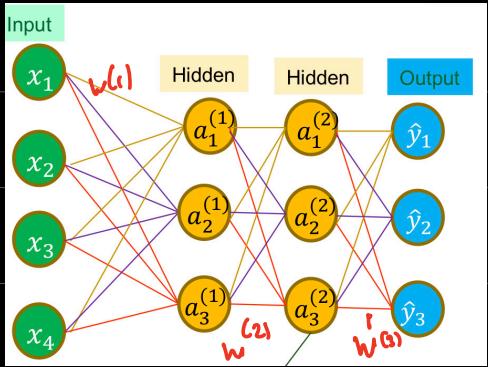


$$\boxed{\frac{\partial J}{\partial w^{(l)}} = \delta^{(l)} a^{(l+1)}}$$

$$\begin{aligned} & \times \xrightarrow{w} \hat{y} \\ & \frac{\partial J}{\partial a_3} = \delta^{(l)} \\ & \frac{\partial J}{\partial w} \leftarrow \underbrace{\delta^{(l)}}_{\hat{y} - y} \quad \underbrace{\frac{\partial J}{\partial a_3}}_{\delta^{(l)}} \\ & z^{(l)} = w^{(l)} a^{(l-1)}; \quad a^{(l)} = g(z^{(l)}) \\ & \delta^{(l)} = e^{(l)} \cdot g'(z^{(l)}) \end{aligned}$$

NLP - Back prop.

Skiping bias



$$\begin{aligned} Z_1 &= w_{11} x_1 + w_{12} x_2 + w_{13} x_3 + w_{14} x_4 \\ Z_2 &= w_{21} x_1 + w_{22} x_2 + w_{23} x_3 + w_{24} x_4 \\ Z_3 &= w_{31} x_1 + w_{32} x_2 + w_{33} x_3 + w_{34} x_4 \end{aligned}$$

$$\begin{bmatrix} Z_1^{(1)} \\ Z_2^{(1)} \\ Z_3^{(1)} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

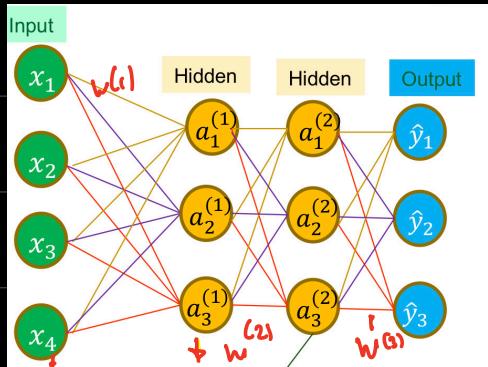
$$Z^{(1)} = W^{(1)} x$$

$$a^{(1)} = g \begin{pmatrix} (z_1^{(1)}) \\ (z_2^{(1)}) \\ (z_3^{(1)}) \end{pmatrix}$$

$$a^{(1)} = g(Z^{(1)})$$

↳ for every component

$$a^{(0)} = x$$



$$Z^{(1)} = \underbrace{W^{(1)} a^{(0)}}_{3 \times 4} \quad a^{(1)} = g(Z^{(1)})$$

$$Z^{(2)} = \underbrace{W^{(2)} a^{(1)}}_{3 \times 3 \times 1}, \quad a^{(2)} = g(Z^{(2)}).$$

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

$$a_i^{(2)}$$



$$\circledcirc \rightarrow \hat{g}_1(a_1^{(2)})$$

$$a_2^{(2)}$$

$$\circledcirc \rightarrow \hat{g}_2(a_2^{(2)})$$

$$a_3^{(2)}$$



$$\circledcirc \rightarrow \hat{g}_3(a_3^{(2)})$$

$$\hat{y} = a^{(3)}$$

$$\frac{\partial J}{\partial z}$$

$$\frac{\partial J}{\partial a} = \epsilon$$

$$\underbrace{\epsilon^{(3)}_{3 \times 1}}_{\text{Can be calculated}} = \frac{\partial J}{\partial a^{(3)}} = \frac{\partial J}{\partial \hat{y}}$$

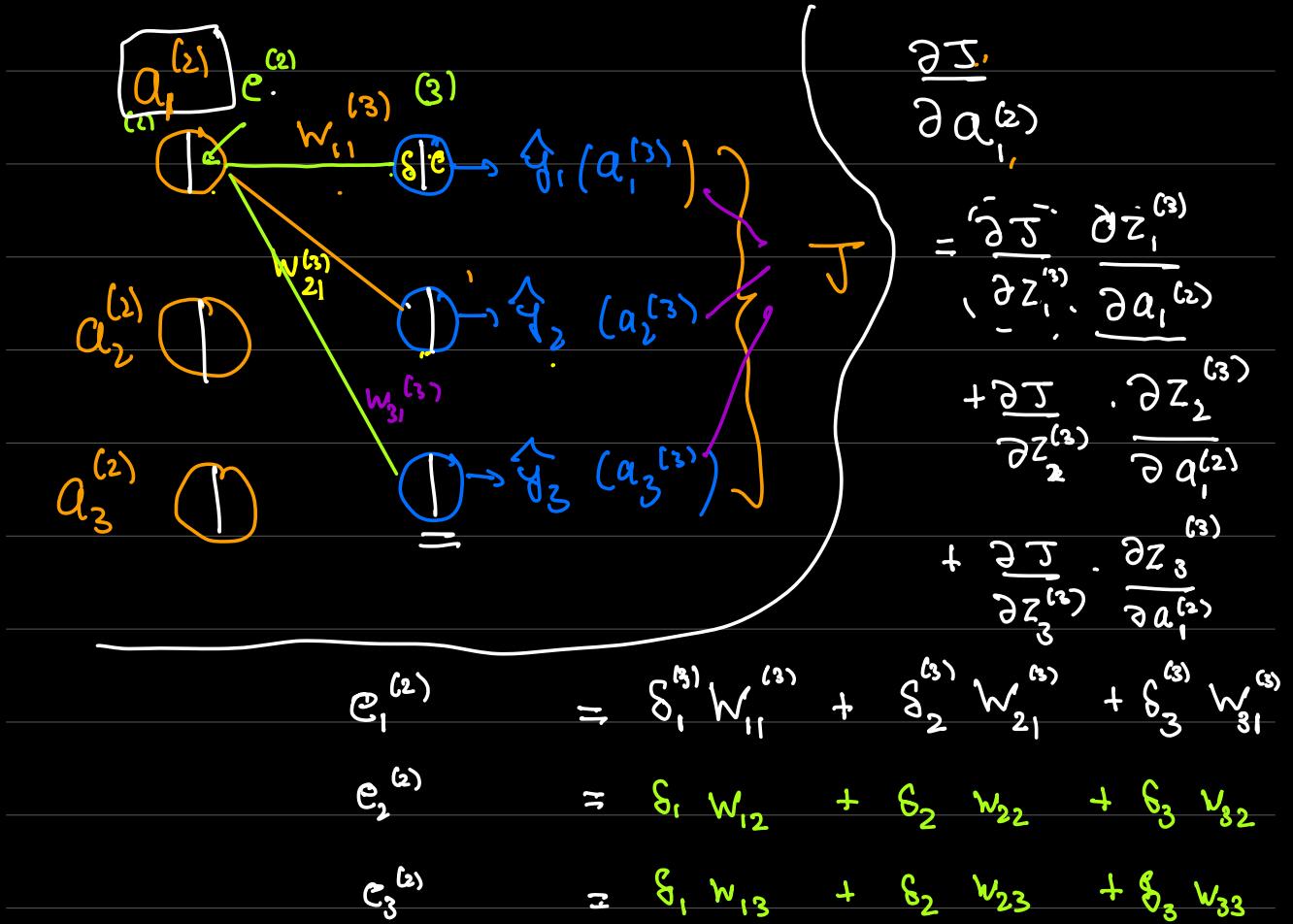
e.g. MSE  
 $\Rightarrow \epsilon = \hat{y} - y$

$$a^{(3)} = g(Z^{(3)}).$$

$$\frac{\partial J}{\partial Z^{(3)}} = \frac{\partial J}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial Z^{(3)}}$$

$$\Rightarrow \frac{\partial a^{(3)}}{\partial Z^{(3)}} = g'(Z^{(3)})$$

$$\underbrace{\epsilon^{(3)}_{3 \times 1}}_{\text{Elementwise}} \circ \underbrace{g'(Z^{(3)})_{3 \times 1}}_{\text{Elementwise}}$$

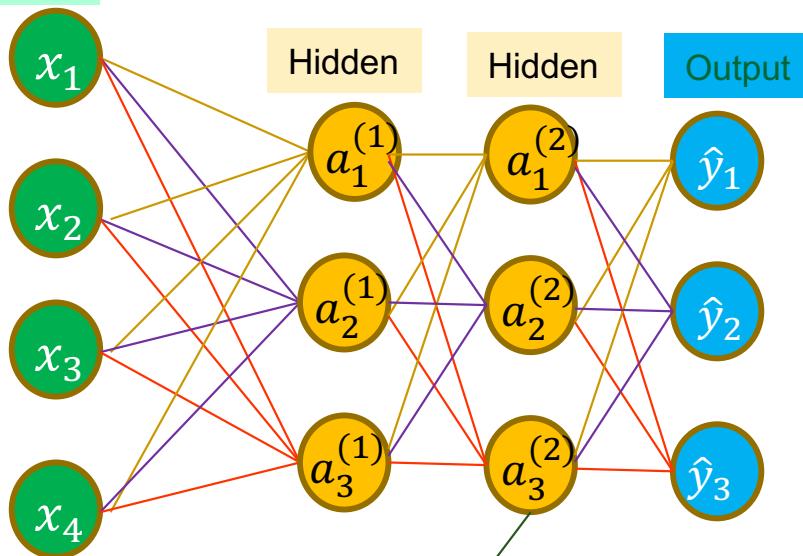


$$\Rightarrow \mathcal{C}^{(2)} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \\ w_{13} & w_{23} & w_{33} \end{bmatrix} \begin{bmatrix} \delta_1 \\ \delta_2 \\ \delta_3 \end{bmatrix}$$

$$\Rightarrow \boxed{\mathcal{C}^{(2)} = W^T \delta^{(3)}}$$

# Forward pass through a deep network

Input



## Main Ideas

Every neuron has a linear and a nonlinear operation

The linear intermediate output is denoted by  $z_i^{(l)}$

The nonlinear final output is denoted by  $a_i^{(l)}$

$$z_i^{(l)} = \sum w_{ij} a_i^{(l-1)} \quad \left. \begin{array}{c} \\ \end{array} \right\} \text{Scalar}$$

Linear

$$\begin{aligned} z_1^{(2)} &= W_{11}a_1^{(1)} + W_{12}a_2^{(1)} + W_{13}a_3^{(1)} \\ z_2^{(2)} &= W_{21}a_1^{(1)} + W_{22}a_2^{(1)} + W_{23}a_3^{(1)} \\ z_3^{(2)} &= W_{31}a_1^{(1)} + W_{32}a_2^{(1)} + W_{33}a_3^{(1)} \end{aligned}$$

Nonlinear

$$\begin{aligned} a_1^{(2)} &= g(z_1^{(2)}) \\ a_2^{(2)} &= g(z_2^{(2)}) \\ a_3^{(2)} &= g(z_3^{(2)}) \end{aligned}$$

Many choices for the nonlinear activation function  $g$  are possible in the intermediate layers

In the output layer  $g$  depends on the task.

$$\begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} = \begin{bmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \\ W_{31} & W_{32} & W_{33} \end{bmatrix} \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \end{bmatrix}$$

$$\mathbf{z}^{(2)} = \mathbf{W}\mathbf{a}^{(1)}$$

$$\begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{bmatrix} = g\left(\begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}\right)$$

$$\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$$

## Forward Pass

$$\text{Step 1. } \mathbf{z}^{(1)} = \mathbf{W}\mathbf{x} \quad \mathbf{a}^{(1)} = g(\mathbf{z}^{(1)})$$

$$\text{Step 2. } \mathbf{z}^{(2)} = \mathbf{W}\mathbf{a}^{(1)} \quad \mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$$

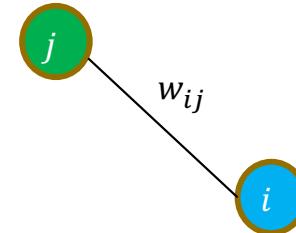
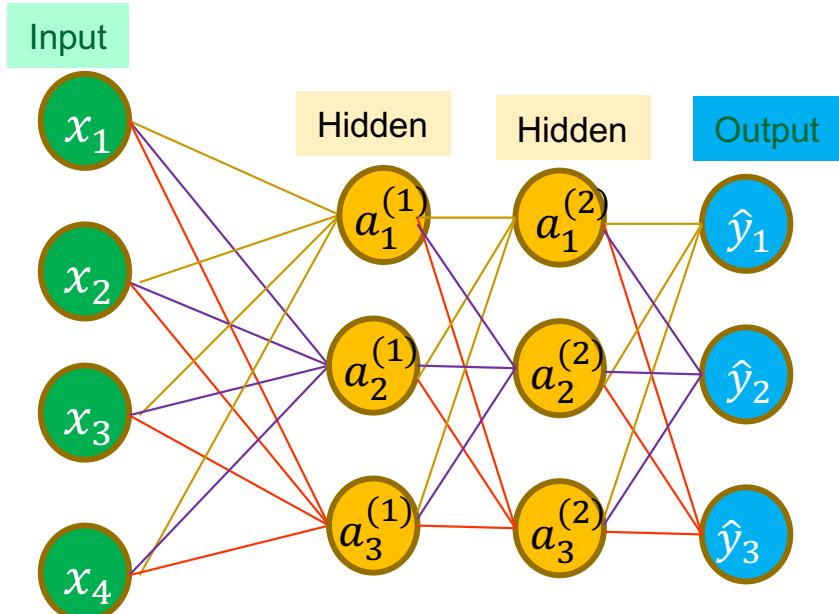
$$\text{Step 3. } \mathbf{z}^{(3)} = \mathbf{W}\mathbf{a}^{(2)} \quad \mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$$

$$\mathbf{a}^{(0)} = \mathbf{x}$$

$$\hat{\mathbf{y}} = \mathbf{a}^{(3)}$$

For  $N_w$  weights, number of multiplications in forward pass  $F \approx O(N_w)$

# FDM or the need for backprop



## Finite Difference

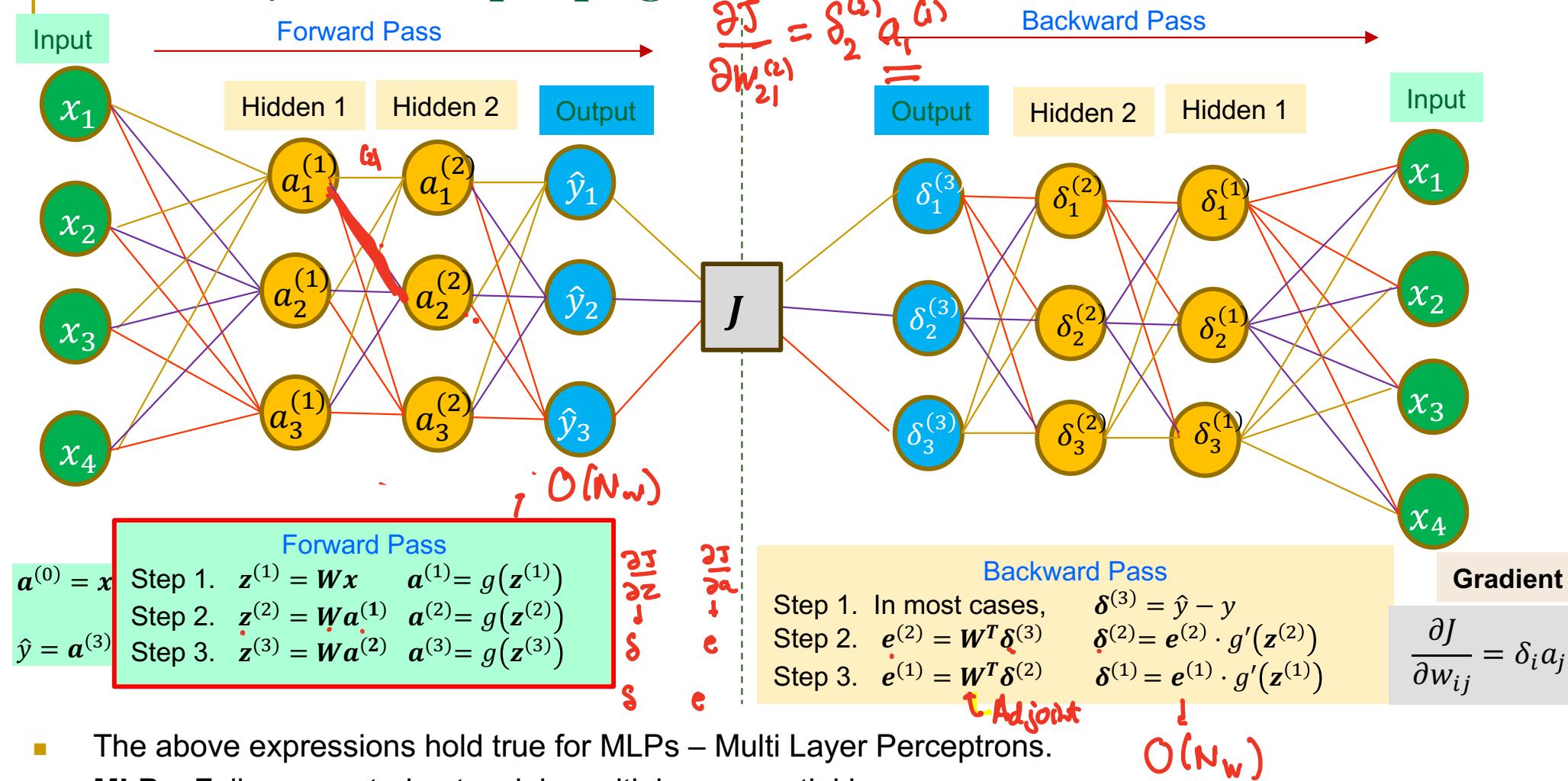
$$\frac{\partial J}{\partial w_{ij}} = \frac{(J(w_{11}, \dots, \textcolor{red}{w_{ij} + \Delta w_{ij}}, \dots) - J(w_{11}, \dots, \textcolor{red}{w_{ij}}, \dots))}{\Delta w_{ij}}$$

- The forward pass gave us  $\hat{y}$  as well as all the  $\hat{a}$  (in the hidden layers)
  - Linear calculations :  $\mathbf{z} = \mathbf{W}\mathbf{a}$  which were matrix multiplications
  - Nonlinear calculations :  $\mathbf{a} = \mathbf{g}(\mathbf{z})$  which are element wise calculations
- What we need for gradient descent is  $w = w - \alpha \frac{\partial J}{\partial w_{ij}}$  for every  $w_{ij}$  in every iteration
- **Finite Difference Method (FDM)** Inefficient way of computing  $\frac{\partial J}{\partial w_{ij}}$ 
  - Remember that gradient descent takes many iterations (epochs) to converge.
  - FDM will take  $(N_w + 1) \times F$  calculations every iteration. (**1 base +  $N_w$  perturbed forward passes**)
  - $N_w$  is the number of weights and  $F$  is the computational expense of one forward pass
  - **$O(N_w^2)$**  Very expensive for big networks. Solution is backprop which is implemented in modern networks

## Forward Pass

- Step 1.  $\mathbf{z}^{(1)} = \mathbf{W}\mathbf{x} \quad \mathbf{a}^{(1)} = g(\mathbf{z}^{(1)})$
- Step 2.  $\mathbf{z}^{(2)} = \mathbf{W}\mathbf{a}^{(1)} \quad \mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$
- Step 3.  $\mathbf{z}^{(3)} = \mathbf{W}\mathbf{a}^{(2)} \quad \mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$

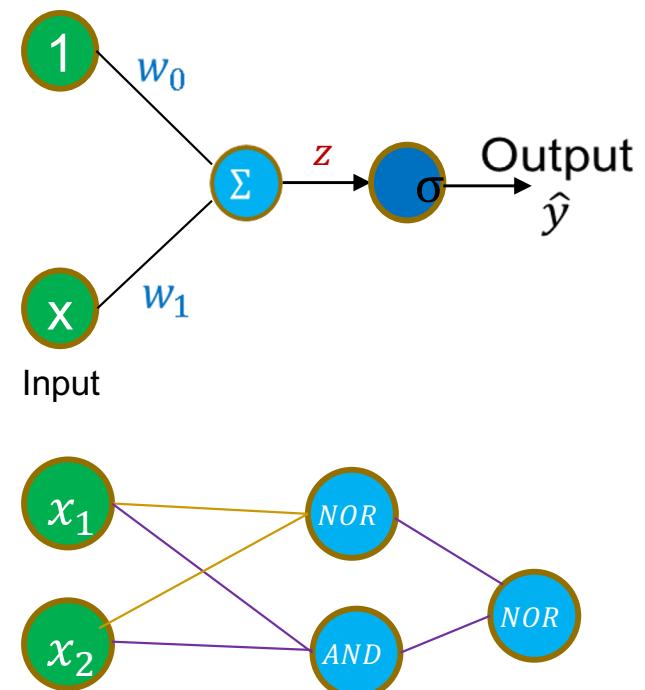
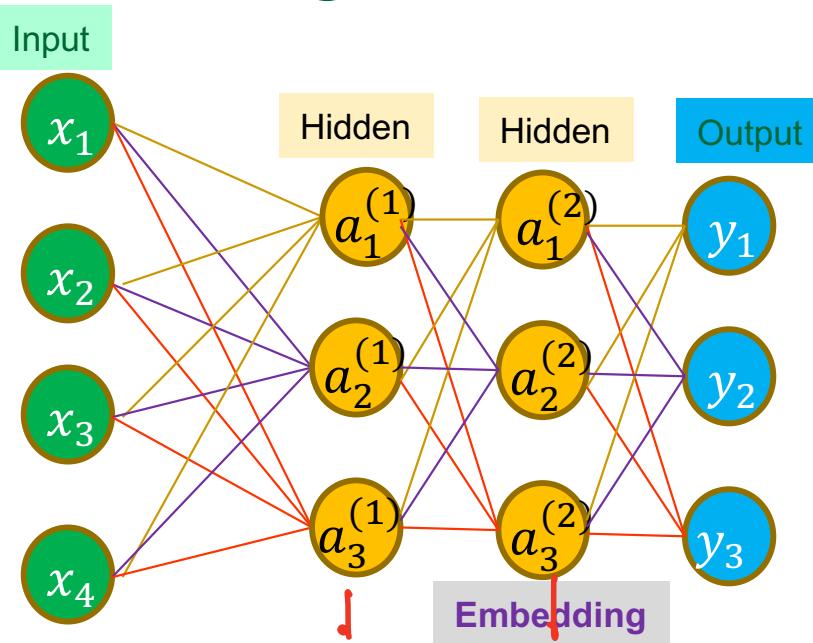
# Summary of backpropagation in MLPs



- The above expressions hold true for MLPs – Multi Layer Perceptrons.
- MLP** – Fully connected network in multiple, sequential layers.
- Notice the symmetry between the forward pass and the backward pass.
- Both involve a sequence of Linear (Matrix multiplication) followed by a nonlinear step ( $g(z)$  or  $g'(z)$ )
- Other Deep Neural Networks need not have such need formulae but the overall symmetry remains

For  $N_w$  weights, number of multiplications in backprop  $F \approx O(N_w)$

# Embeddings – Or how deep neural networks work

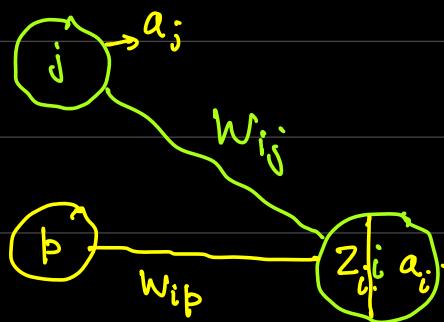


- Every “deep” network has, in essence, **a prefix** to a final linear, logistic, or softmax layer
  - That is, a linear regression or a linear classification problem
- The task of all the prefix (hidden) layers is to transform the problem into a representation which will work linearly.
- That is, **deep networks learn representations** that work simply with final layers
- The final representation before the output layer is called the **Embedding**
- **Embeddings** are like zipped, vector versions of the input
  - Of great practical importance. OpenAI “sells” embeddings for ChatGPT currently

$$z^{(3)} = \underline{w} \underline{a}^{(2)} ; \quad a_j^{(2)} = w^{-1} z^{(3)}$$

$$\frac{\partial J}{\partial a^{(2)}} = w^{-1} \frac{\partial J}{\partial z}$$

$$\frac{\partial J}{\partial z}, \quad \frac{\partial J}{\partial a} = \underline{w^T} \underline{s}^{(3)}$$



$$\frac{\partial J}{\partial w_{ij}} = \underbrace{\frac{\partial J}{\partial z_i}}_{\delta_i} \underbrace{\frac{\partial z_i}{\partial w_{ij}}}_{a_j}.$$

$$\boxed{\frac{\partial J}{\partial w_{ij}} = \delta_i a_j}$$

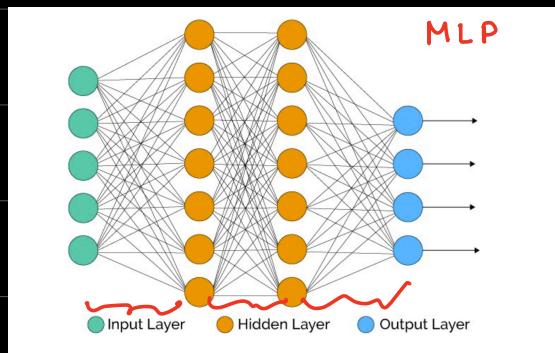
"DELTA" RULE

$$Z_i = w_{ij} a_j + \underbrace{w_{ip} a_p}_0$$

$$\frac{\partial Z_i}{\partial w_{ij}} = a_j + 0$$

## RECAP

### Expressions in the general, MLP case



$$\text{Final Layer : } e = \frac{\partial J}{\partial y}$$

Hidden  
Layers

$$e^{(k)} = W^T s^{(k+1)}$$

$$s^{(k)} = g'(z) \odot e^{(k)}$$

$$\frac{\partial J}{\partial w} \equiv \Delta w$$

$$\frac{\partial J}{\partial a} \equiv e$$

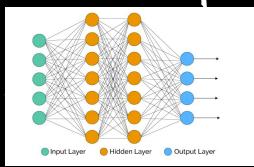
Backprop  
Gradient

$$\frac{\partial J}{\partial z} \equiv \delta$$

### Full Algorithm for Backprop weight update

1. Initialize Weights Randomly

2. Perform Forward Prop till final layer → For every data pt  
 Linear:  $Z^{(k)} = W^{(k)} a^{(k-1)}$       Nonlinear:  $a^{(k)} = g(Z^{(k)})$



3. Calculate  $e = \frac{\partial J}{\partial y}$  at the output of the final layer

4. Perform Backprop till first layer

$$\text{Nonlinear: } \delta^{(k)} = g'(Z^{(k)}) \odot e^{(k)}$$

$$\frac{\partial J}{\partial z} \quad \frac{\partial J}{\partial a}$$

$$\text{Linear: } e^{(k)} = W^{(k+1)} \delta^{(k+1)}$$

$$\frac{\partial J}{\partial w} \quad \frac{\partial J}{\partial v}$$

5. Use the "Delta" Rule to calculate the weight update

$$\frac{\partial J}{\partial w_{ij}} = \delta_i a_j \equiv \Delta w_{ij}$$

6. Use Gradient Descent to update  $w_{ij}$

$$w_{ij} = w_{ij} - \alpha \Delta w_{ij}$$

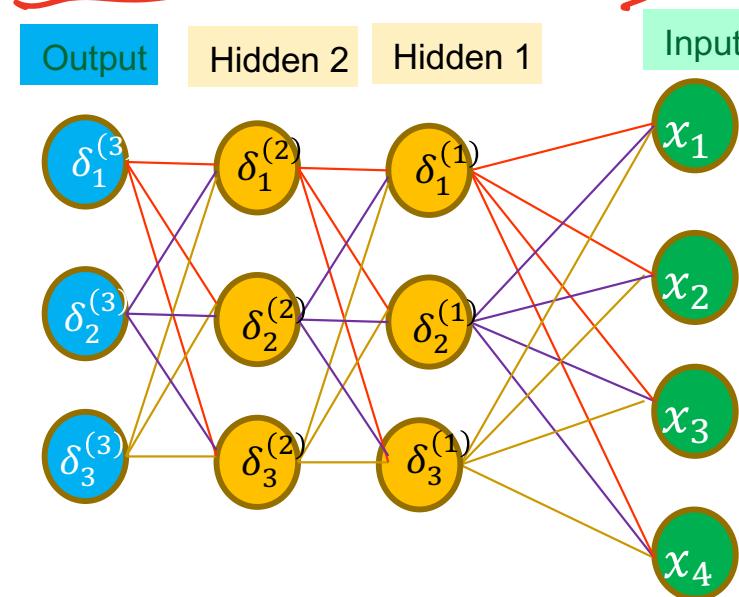
7. Repeat Steps 2-4 for all data pts → One Epoch

8. Repeat Steps 2-5 till convergence → Multiple Epochs

Problems in Deep Networks : Vanishing Gradient, Exploding Gradient, Overfitting  
 Activation Function      RNNs      Regularization  
 Dropout

# Practical Consideration – Activation Functions, Nonlinearity

- When people tried to train deep networks with  $g = \sigma$ , the weights did not update
- Reason – In the backprop formula, we have  $g'(z)$  used recursively
- This means that if  $g'(z) < 1$ , the gradient rapidly decays over the layers – **Vanishing Gradient Problem**
- Solution is to use other nonlinearities  $g(z) = \tanh$  and  $g(z) = \text{ReLU}(z)$

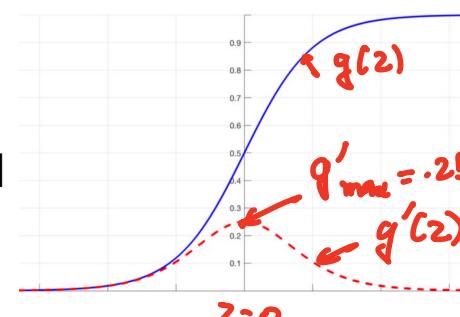
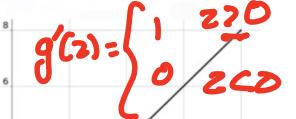
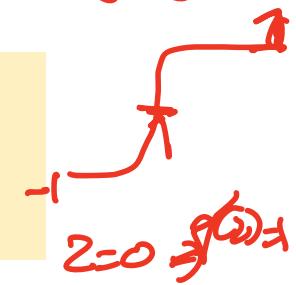


$$g(z) = \tanh(z)$$

$$\frac{e^z - e^{-z}}{e^z + e^{-z}}$$

**Backward Pass**

Step 1. In most cases,	$\delta^{(3)} = \hat{y} - y$
Step 2. $e^{(2)} = W^T \delta^{(3)}$	$\delta^{(2)} = e^{(2)} \cdot g'(z^{(2)})$
Step 3. $e^{(1)} = W^T \delta^{(2)}$	$\delta^{(1)} = e^{(1)} \cdot g'(z^{(1)})$



$$g(z) = \sigma; \quad g'(z) = \sigma(1-\sigma)$$

$$\cdot 5 * (1-0.5) = 0.25$$

