

# Discriminative Models II

Saumitra Yadav

# Last Lecture

Problem: Classify wolf and Brynden Rivers (as a raven in dream) (POS Tagging)

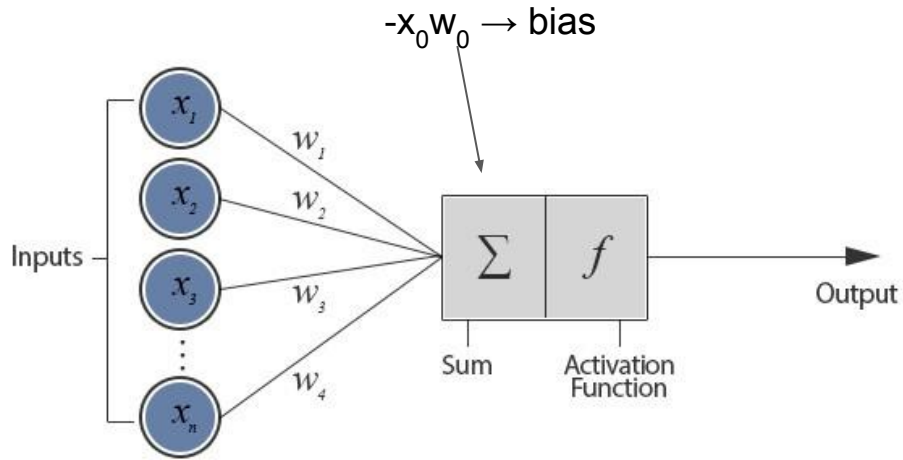
Generative Models (e.g: naive bayes):

- What does wolf and Brynden Rivers (as a raven in dream) look like? ( Generating Words given POS Tags )

Discriminative Models (e.g: logistic regression, neural networks):

- Find features ( one lacks wings, a third eye and another one doesn't stand on all limbs and is present in a dream) to discriminate between inputs. (what kind of distinctions do nouns have when compared with verb- distinctions like company they keep etc.)

# Artificial Neuron



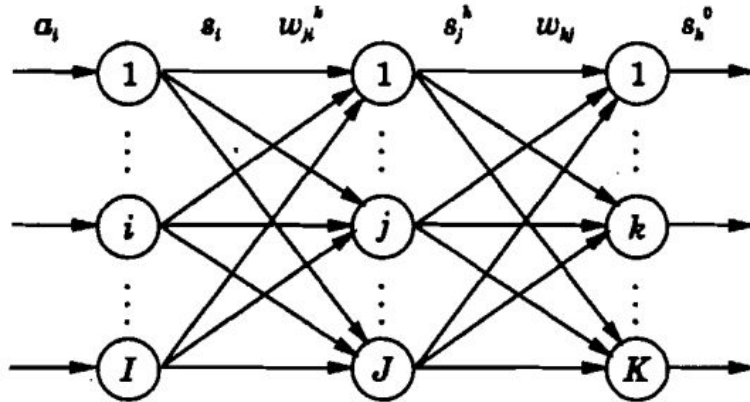
Taken From [1]

# Perceptron Learning Algorithm (Frank Rosenblatt)

- Working of neural network governed by neural dynamics.
  - a. Short term memory - Activation Dynamics - (Sum :  $\sum w_i x_i$ )
  - b. Long term memory - Synaptic Dynamics - (  $W$  - containing all weights )
- Learning Laws are implementation models of Synaptic Dynamics
- Algorithm
  - a. Initialise  $W$  to small random values
  - b. Present a vector  $x$  as input and calculate the output  $y$ .
  - c. Updates the weights according to
    - $w_j(t+1) = w_j(t) + \eta(y' - y)x$
    - Where
      - $y'$  is desired output
      - $t$  is the iteration number
      - $\eta$  is the step size in range 0.0 and 1.0
  - d. Repeat b and c until one of the following conditions is met:
    - Iteration error ( square of summation of all errors ) is less than a threshold value
    - $t$  has reached a predetermined number for completion.

# Multilayer Perceptron Network

- Classify non-linear separable data



$a_i$ :  $i$ th unit of  $I$  dimensional input

$s_i$ : output of activation of  $i$ th unit in input layer

$w_{ji}^h$ : weight from  $i$ th unit to  $j$ th unit; where  $j$ th unit is in hidden layer

$s_j^h$ : output of activation of  $j$ th unit in hidden layer

$w_{kj}$ : weight from  $j$ th hidden unit to  $k$ th unit in output layer

$s_k^o$ : output of activation of  $k$ th unit in output layer

Taken from [2]

# Backpropagation Learning Algorithm

- Network is a function mapping from input to output space  $\mathbf{R}^I : \mathbf{R}^K$
- Using gradient search along error surface for optimum sets of weight
- Given input  $\mathbf{a}(\mathbf{m})$ , calculated output  $\mathbf{b}'(\mathbf{m})$  and  $\mathbf{b}(\mathbf{m})$  being actual output in  $\mathbf{m}_{th}$  iteration
- Let  $\mathbf{E}(\mathbf{m})$  be error, then gradient descent along error surface so as to change weight connection units  $j$  and  $i$ .

$$\Delta w_{ij}(m) = -\eta \frac{\partial E(m)}{\partial w_{ij}} \quad \text{Where } \eta > 0$$

- Accordingly weight update will be as follows

$$w_{ij}(m+1) = w_{ij}(m) + \Delta w_{ij}(m)$$

# Backpropagation Learning Algorithm Contd.

- Error (mean squared error) for  $m^{\text{th}}$  step is,

$$\begin{aligned} E(m) &= \frac{1}{2} \sum_{k=1}^K [b_k(m) - b'_k(m)]^2 = \frac{1}{2} \sum_{k=1}^K [b_k(m) - s_k^o]^2 \\ &= \frac{1}{2} \sum_{k=1}^K [b_k(m) - f_k^o(x_k^o)]^2, \end{aligned}$$

- Where

$$x_k^o = \sum_{j=1}^J w_{kj} s_j^h$$

$$s_j^h = f_j^h(x_j^h)$$

$$x_j^h = \sum_{i=1}^I w_{ji}^h s_i$$

$$s_i = x_i = a_i(m).$$

# Backpropagation Learning Algorithm Contd.

For weights leading to output unit

$$\Delta w_{kj}(m) = -\eta \frac{\partial E(m)}{\partial w_{kj}}$$

$$\begin{aligned}\frac{\partial E(m)}{\partial w_{kj}} &= \frac{1}{2} \frac{\partial}{\partial w_{kj}} \left[ b_k - f_k^o \left( \sum_{j=1}^J w_{kj} s_j^h \right) \right]^2 \\ &= -(b_k - f_k^o) \dot{f}_k^o s_j^h \\ &= -\delta_k^o s_j^h,\end{aligned}$$

Where,

$$\delta_k^o = (b_k - f_k^o) \dot{f}_k^o.$$



# Backpropagation Learning Algorithm Contd.

Therefore  $\Delta w_{kj}(m) = \eta \delta_k^o s_j^h$

And,

$$\begin{aligned} w_{kj}(m+1) &= w_{kj}(m) + \Delta w_{kj}(m) \\ &= w_{kj}(m) + \eta \delta_k^o s_j^h. \end{aligned}$$



Update is  
learning\_rate x  
error propagated  
back by  $k^{\text{th}}$  unit x  
activation output at  
j

Gradient for weights leading to units in hidden layer is,

$$\Delta w_{ji}^h(m) = -\eta \frac{\partial E(m)}{\partial w_{ji}^h}$$

# Backpropagation Learning Algorithm Contd.

Which gives us,

$$\begin{aligned}\frac{\partial E(m)}{\partial w_{ji}^h} &= -\sum_{k=1}^K (b_k - f_k^o) \frac{\partial f_k^o \left( \sum_{j=1}^J w_{kj} s_j^h \right)}{\partial w_{ji}^h} \\ &= -\sum_{k=1}^K (b_k - f_k^o) f_k^o w_{kj} \frac{\partial f_k^h}{\partial w_{ji}^h},\end{aligned}$$

And since  $s_j^h = f_j^h(x_j^h)$

$$\frac{\partial s_j^h}{\partial w_{ji}^h} = f_j^h \frac{\partial x_j^h}{\partial w_{ji}^h}$$

# Backpropagation Learning Algorithm Contd.

And Since  $x_j^h = \sum_{i=1}^I w_{ji}^h s_i$  so,  $\frac{\partial x_j^h}{\partial w_{ji}^h} = s_i$

Hence,

$$\begin{aligned} \frac{\partial E(m)}{\partial w_{ji}^h} &= - \sum_{k=1}^K (b_k - f_k^o) \dot{f}_k^o w_{kj} \dot{f}_j^h s_i \\ &= - \delta_j^h s_i \end{aligned} \quad \text{Where, } \delta_j^h = \dot{f}_j^h \sum_{k=1}^K w_{kj} \delta_k^o$$

# Backpropagation Learning Algorithm Contd.

Hence,

$$\Delta w_{ji}^h(m) = \eta \delta_j^h s_i = \eta \delta_j^h a_i(m)$$

And since  $s_i = x_i = a_i(m)$  therefore,

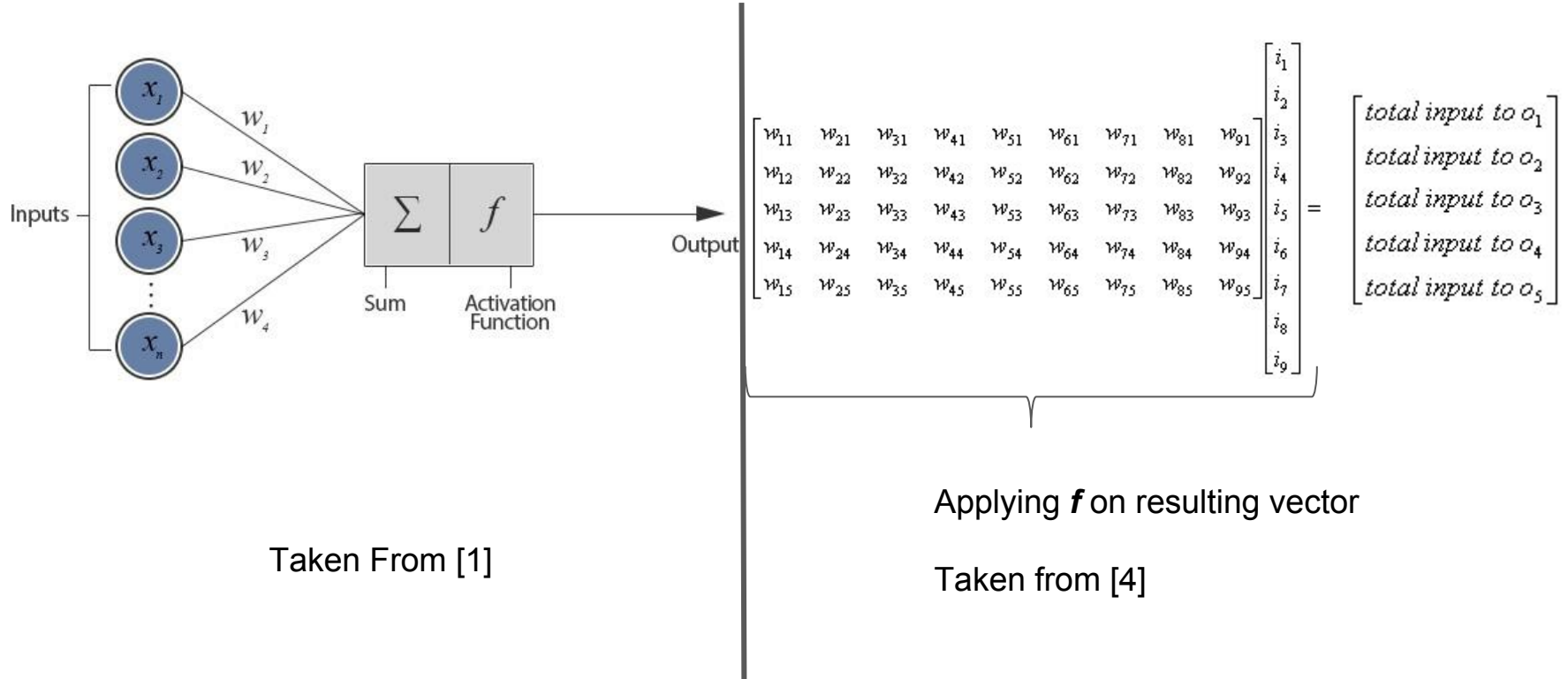
$$\begin{aligned} w_{ji}^h(m+1) &= w_{ji}^h(m) + \Delta w_{ji}^h(m) \\ &= w_{ji}^h(m) + \eta \delta_j^h a_i(m) \end{aligned}$$

Where,

$$\delta_j^h = f_j^h \sum_{k=1}^K w_{kj} \delta_k^o, \text{ denotes the error propagated back.}$$

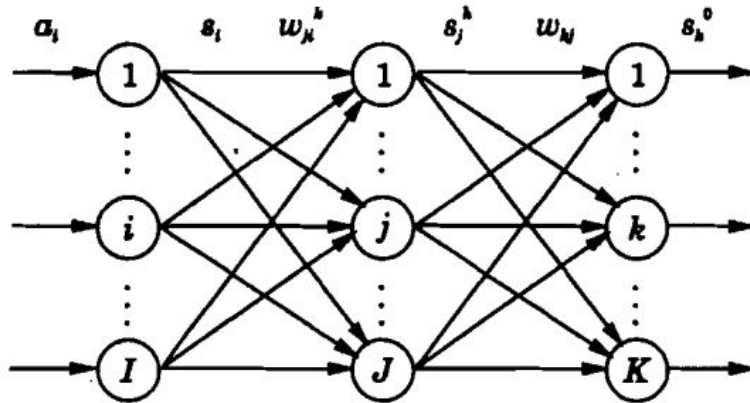
**Update** is learning\_rate  $\times$  error propagated back from  $j^{\text{th}}$  unit ( which consist of sum of all error propagated back to  $j^{\text{th}}$  unit from all units in layer **K** )  $\times$   $i^{\text{th}}$  output of input layer

# Artificial Neuron - Matrix Perspective



# Multilayer Perceptron Network - Matrix Perspective

- Classify non-linear separable data



$a_i$ :  $i$ th unit of  $I$  dimensional input  
 $s_i$ : output of activation of  $i$ th unit in input layer  
 $w_{ij}^h$ : weight from  $i$ th unit to  $j$ th unit; where  $j$ th unit is in hidden layer  
 $s_j^h$ : output of activation of  $j$ th unit in hidden layer  
 $w_{jk}^o$ : weight from  $j$ th hidden unit to  $k$ th unit in output layer  
 $s_k^o$ : output of activation of  $k$ th unit in output layer

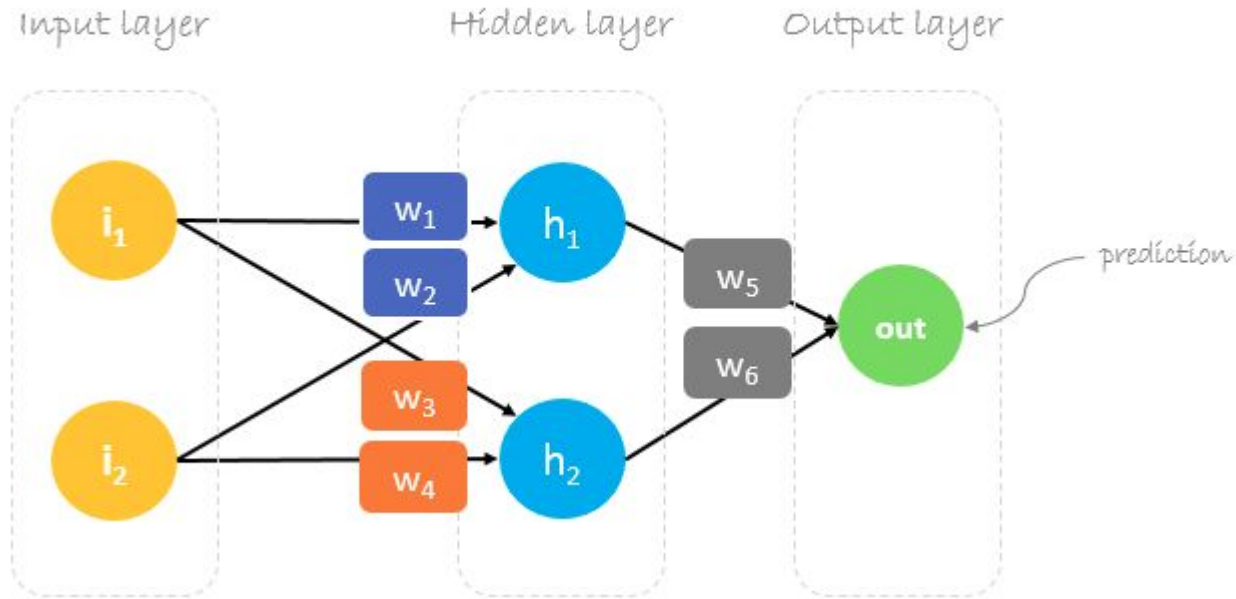
$$\text{output} = f^o(W^o x f^h(W^h x A))$$

Where,

$A$ : input vector/matrix of size  $n \times v$ ,  $n$  is number of dimension and  $v$  is 1/number of samples.  
 $W^h$ : Weight Matrix between hidden and input layer  
 $f^h$ : activation function of hidden layer  
 $W^o$ : Weight Matrix between hidden and output layer  
 $f^o$ : activation function of output layer at each cell

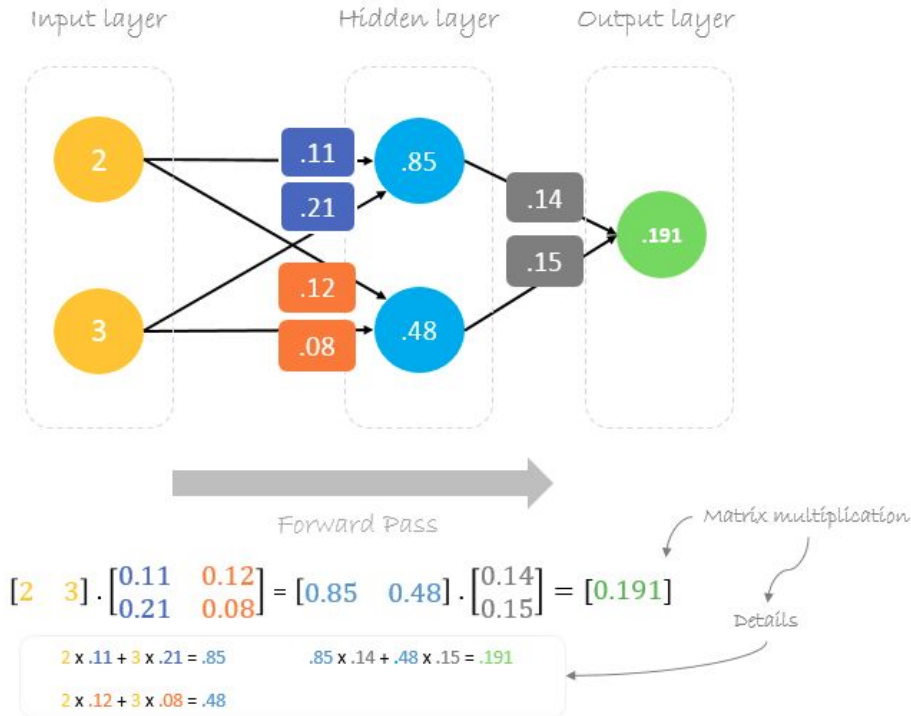
Taken from [2]

# Matrix Perspective naive example



Taken from [3]

# Matrix Perspective naive example contd.



Actual Output : 1  
Error =  $(0.191 - 1)^2/2$   
= 0.327

Derivative of Error with respect to weight

Old weight

New weight

Learning rate

$$*W_x = W_x - a \left( \frac{\partial \text{Error}}{\partial W_x} \right)$$



# Matrix Perspective naive example contd.

$$\frac{\partial \text{Error}}{\partial W_6} = \frac{\partial \text{Error}}{\partial \text{prediction}} * \frac{\partial \text{prediction}}{\partial W_6} \quad \leftarrow \text{chain rule}$$

$$\text{Error} = \frac{1}{2}(\text{prediction} - \text{actual})^2$$

$$\frac{\partial \text{Error}}{\partial W_6} = \frac{\frac{1}{2}(\text{prediction} - \text{actual})^2}{\partial \text{prediction}} * \frac{\partial (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6}{\partial W_6}$$

$$\text{prediction} = (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6$$

$$\frac{\partial \text{Error}}{\partial W_6} = 2 * \frac{1}{2}(\text{prediction} - \text{actual}) \frac{\partial (\text{prediction} - \text{actual})}{\partial \text{prediction}} * (i_1 w_3 + i_2 w_4)$$

$$h_2 = i_1 w_3 + i_2 w_4$$

$$\frac{\partial \text{Error}}{\partial W_6} = (\text{prediction} - \text{actual}) * (h_2)$$

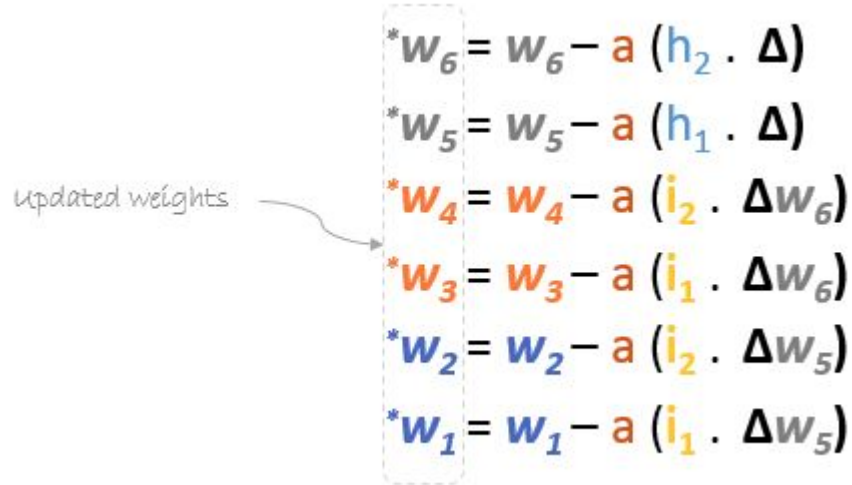
$$\Delta = \text{prediction} - \text{actual}$$

delta

$$\frac{\partial \text{Error}}{\partial W_6} = \Delta h_2$$

# Matrix Perspective naive example contd.

updated weights


$$\begin{aligned} *w_6 &= w_6 - a (h_2 \cdot \Delta) \\ *w_5 &= w_5 - a (h_1 \cdot \Delta) \\ *w_4 &= w_4 - a (i_2 \cdot \Delta w_6) \\ *w_3 &= w_3 - a (i_1 \cdot \Delta w_6) \\ *w_2 &= w_2 - a (i_2 \cdot \Delta w_5) \\ *w_1 &= w_1 - a (i_1 \cdot \Delta w_5) \end{aligned}$$

# Matrix Perspective naive example contd.

$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - a \Delta \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - \begin{bmatrix} a h_1 \Delta \\ a h_2 \Delta \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - a \Delta \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} \cdot [w_5 \quad w_6] = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - \begin{bmatrix} a i_1 \Delta w_5 & a i_1 \Delta w_6 \\ a i_2 \Delta w_5 & a i_2 \Delta w_6 \end{bmatrix}$$

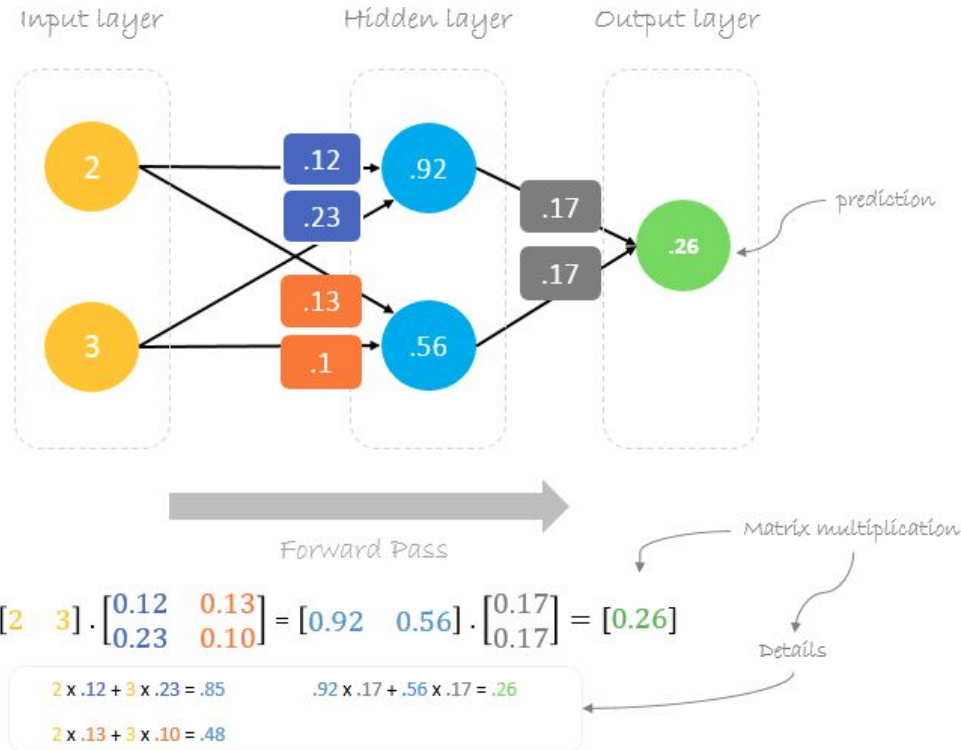
$$\Delta = 0.191 - 1 = -0.809 \quad \leftarrow \text{Delta} = \text{prediction} - \text{actual}$$

$$a = 0.05 \quad \leftarrow \text{Learning rate, we smartly guess this number}$$

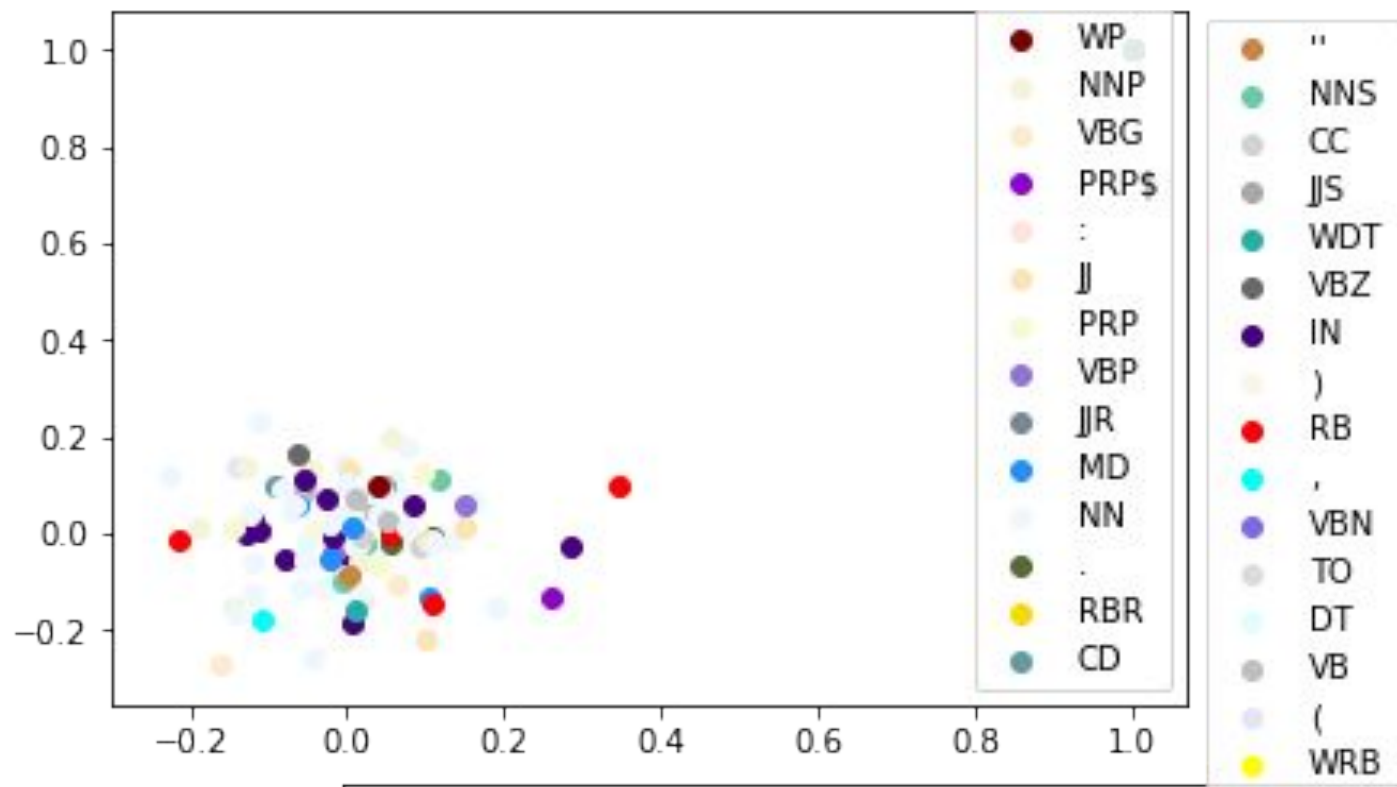
$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} - 0.05(-0.809) \begin{bmatrix} 0.85 \\ 0.48 \end{bmatrix} = \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} - \begin{bmatrix} -0.034 \\ -0.019 \end{bmatrix} = \begin{bmatrix} 0.17 \\ 0.17 \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} .11 & .12 \\ .21 & .08 \end{bmatrix} - 0.05(-0.809) \begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot [0.14 \quad 0.15] = \begin{bmatrix} .11 & .12 \\ .21 & .08 \end{bmatrix} - \begin{bmatrix} -0.011 & -0.012 \\ -0.017 & -0.018 \end{bmatrix} = \begin{bmatrix} .12 & .13 \\ .23 & .10 \end{bmatrix}$$

# Matrix Perspective naive example contd.



# One Application



# References for Text and Images

1. <https://medium.com/technologymadeeasy/for-dummies-the-introduction-to-neural-networks-we-all-need-c50f6012d5eb>
2. Artificial Neural Network by B.Yegnanarayana
3. <http://hmkcode.github.io/ai/backpropagation-step-by-step/>
4. [http://www.sharetechnote.com/html/EngMath\\_Matrix\\_NeuralNetwork.html](http://www.sharetechnote.com/html/EngMath_Matrix_NeuralNetwork.html)

# If you want to further read about MLP (specifically backprop)

1. <https://www.analyticsvidhya.com/blog/2017/05/neural-network-from-scratch-in-python-and-r/>
2. <https://medium.com/datathings/neural-networks-and-backpropagation-explained-in-a-simple-way-f540a3611f5e>
3. <https://medium.com/@14prakash/back-propagation-is-very-simple-who-made-it-complicated-97b794c97e5c>