

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

# Neural Networks & Deep Learning

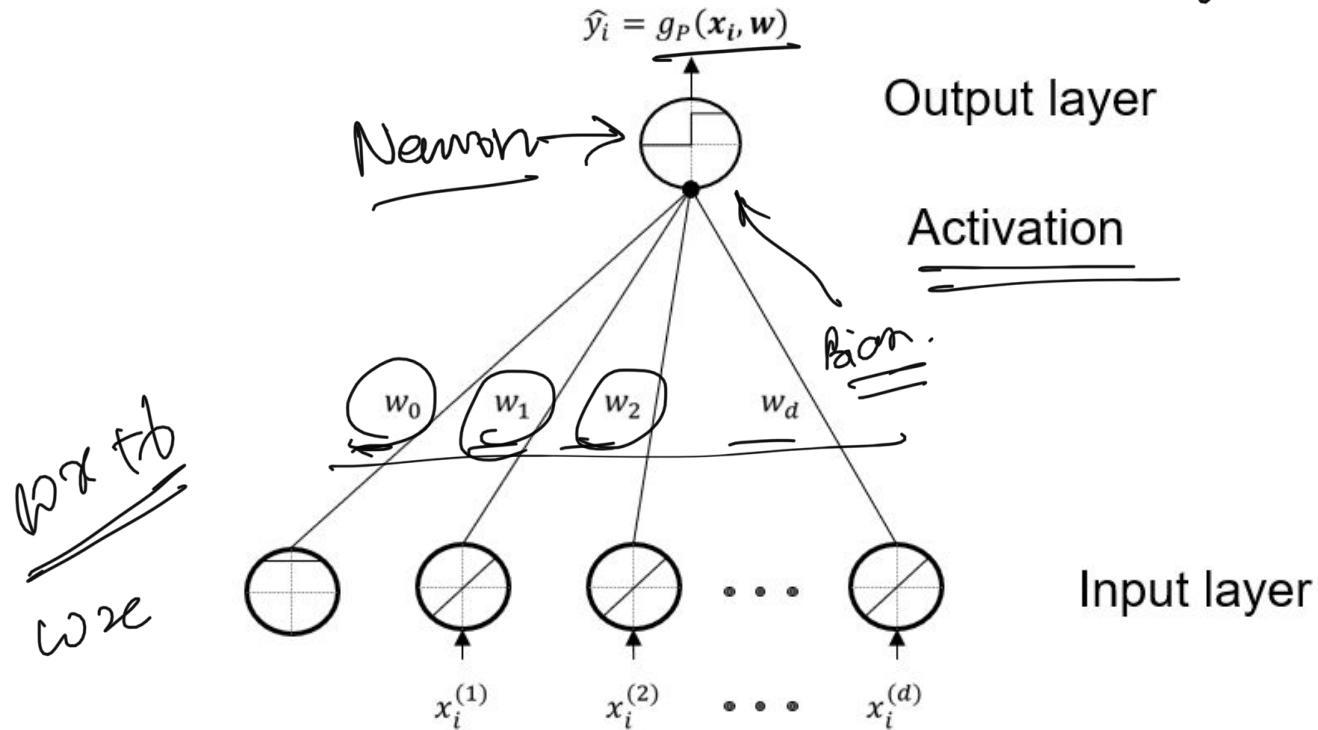


# Table of contents

1. Perceptron Learning
2. XOR Problem using MLP
3. Back-propagation

# What is a Perceptron?

Linear binary classifier.



# Learning in Perceptron

Iterate for  $t$  until a stop criterion is met

{

for each sample  $x_i$  with label  $y_i$ :

{

compute the output  $\hat{y}_i$  of the network

estimate the error of the network  $e(w(t)) = y_i - \hat{y}_i$

update the weight  $w(t+1) = w(t) + e(w(t))x_i$

}

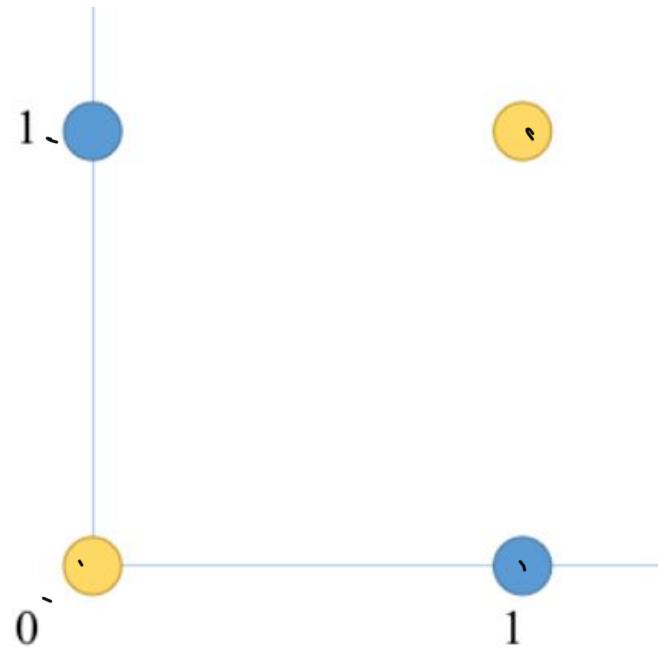
$t++$

}

Max iterations  
→ does not change  
 $\|w_{old} - w_{new}\| < \epsilon$   
ground truth (labels).



# XOR problem



# XOR problem - Using MLP (Multi-layer perceptron) .

$$w_{t+1} = w_t - \eta \nabla_w J$$

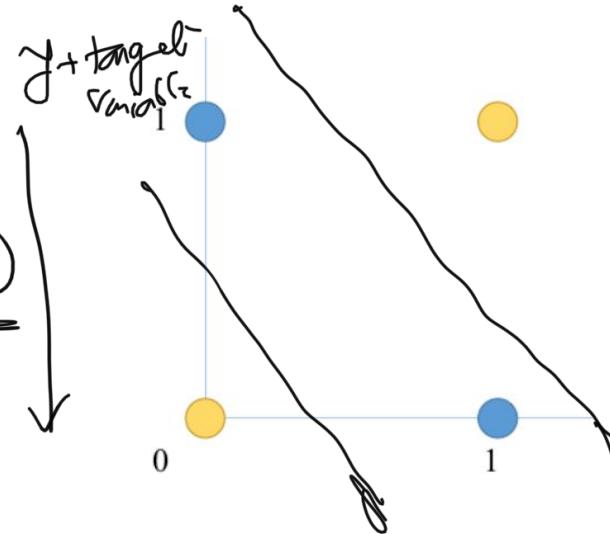
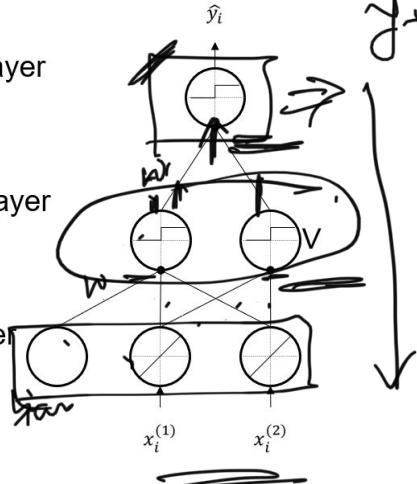
Output layer

Hidden layer

$$z_1 = w_1 x + b_1$$

Input layer

$$z_2 = w_2 x + b_2$$



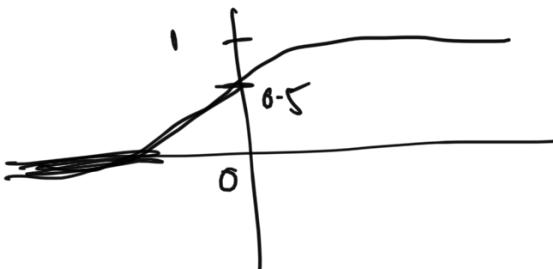
- How do we learn weights? Can perceptron learning be applied?

~ Gradient descent

– Differentiable.

– Backpropagation .

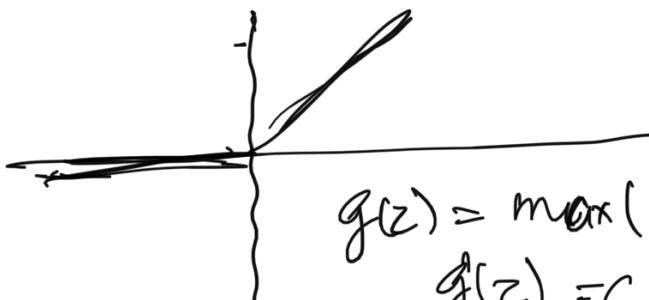
# Activation Functions



Sigmoid

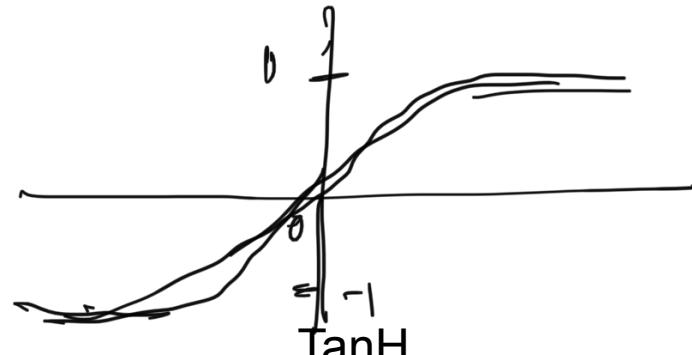
$$g(z) = \frac{1}{1+e^{-z}}$$

$$g'(z) = g(z)(1-g(z))$$



ReLU  
(Rectified Linear Unit)

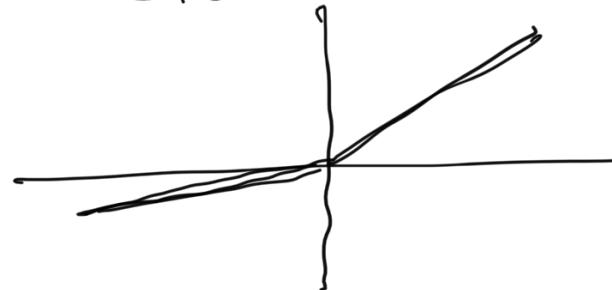
$$g(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$



TanH

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

$$g(z) = 1 - (g(z))^2$$



Leaky ReLU

$$g(z) = \max(\epsilon z, z)$$

# Why non-linear activation functions are important?



$$\underbrace{w_1x + b_1}_{w_2(\sigma(w_1x + b_1))} \rightarrow \underbrace{w_2(w_1x + b_1) + b_2}_{w'(w_1x + b_1)}$$

- Bring non-linearity to the network.

# Handling Multiple Classes

$$\begin{bmatrix} \underline{0.05} & \underline{0.05} & \underline{0.9} \end{bmatrix}$$

- No. of output neurons - No. of classes.
- Activation function - Softmax activation -
  - Bring the values b/w 0 & 1.
  - argmax of <sup>softmax</sup> probabilities
- Error Function-

$C \rightarrow$  no. of classes.

$$\frac{e^{z_i}}{\sum_{j=1}^C e^{z_j}}$$

Normalisation Constant

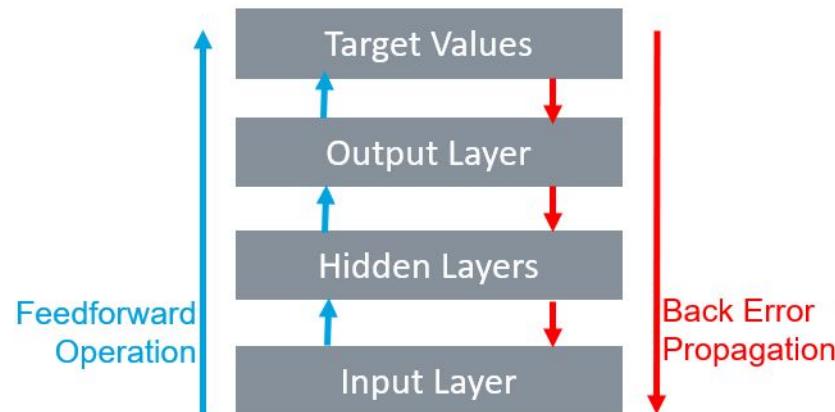
Cross-entropy (Binary cross entropy)

$$L = - \sum_{i=1}^n \sum_{j=1}^C (\gamma_i) \log p(y_i=j | x_i) \rightarrow \begin{cases} 1 & \gamma_i=j \\ 0 & \text{otherwise} \end{cases}$$

softmax probability.

# MLP Learning - Back Propagation

- Key idea - Properly distribute error computed from output layer back to earlier layers to allow their weights to be updated in a way that reduce the error.

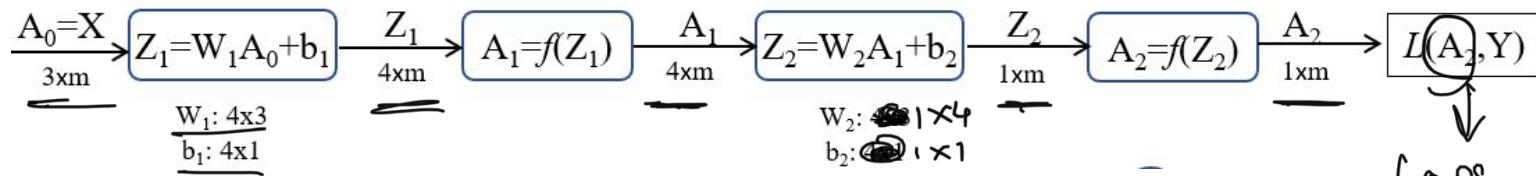
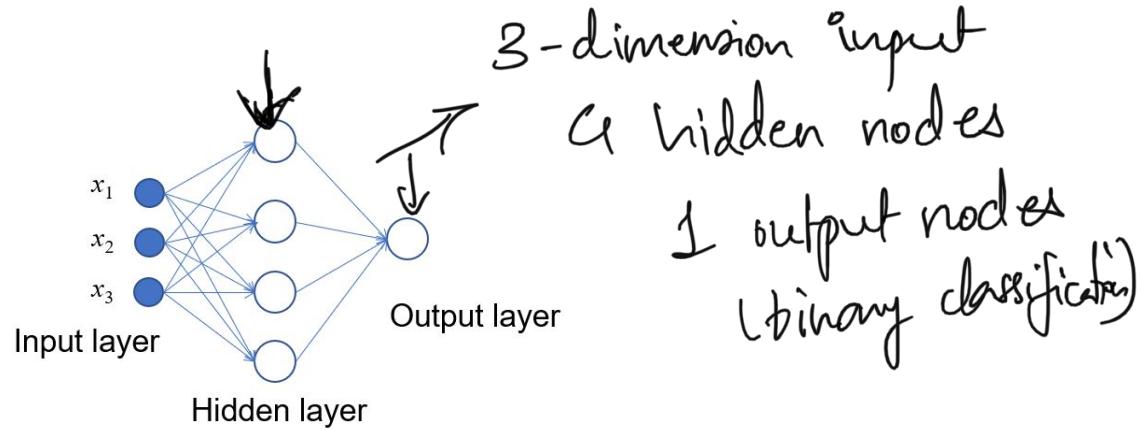


# MLP Learning - Back Propagation

$4 \times 3$        $(4 \times 1)$        $4 \times 3$   
 $\xrightarrow{=}$        $\xrightarrow{\text{Broadcasting}}$

Feed-Forward Phase-

$m \rightarrow$  total no. of samples -  
activation fun<sup>n</sup>  $\rightarrow$  sigmoid.



$$Z_1 = W_1 A_0 + b_1$$

$\downarrow$   
 $4 \times m$   
 $2 \times 3$        $3 \times m$        $4 \times 1$

$$Z_2 = f_2 A_1 + b_2$$

$\downarrow$   
 $1 \times m$   
 $1 \times 4$        $4 \times m$        $1 \times 1$

Loss function

# MLP Learning - Back Propagation

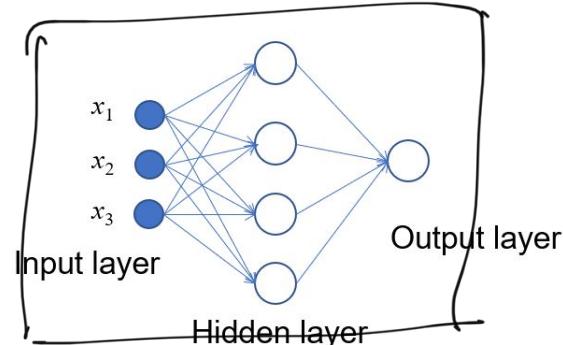
Back propagation phase-

$$w_{t+1} = w_t - \eta \nabla_w L$$

Parameters -

$$w_1, b_1, w_2, b_2$$

L ← Loss function



$$A_0 = X \xrightarrow[3xm]{Z_1 = W_1 A_0 + b_1} \xrightarrow[4xm]{A_1 = f(Z_1)} \xrightarrow[4xm]{Z_2 = W_2 A_1 + b_2} \xrightarrow[1xm]{A_2 = f(Z_2)} \xrightarrow[1xm]{L(A_2, Y)}$$

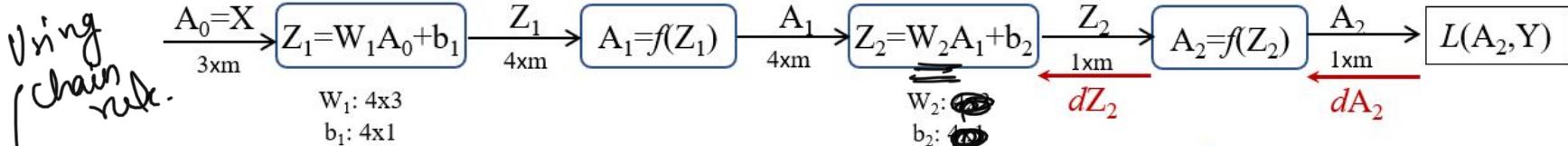
$$W_1: 4 \times 3 \\ b_1: 4 \times 1$$

$$W_2: 1 \times 4 \\ b_2: 1 \times 1$$

$$\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial b_1}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial b_2} \Rightarrow \text{Chain rule. } \frac{\partial L}{\partial A_2}$$

Activation  $\rightarrow$  Sigmoid, Output + Binary loss  $\rightarrow$  BCE

## MLP Learning - Back Propagation



$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial Z_2} \cdot \frac{\partial Z_2}{\partial w_2}, \quad \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial Z_2} \cdot \frac{\partial Z_2}{\partial b_2}, \quad \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial Z_2} \cdot \frac{\partial Z_2}{\partial A_1} \cdot \frac{\partial A_1}{\partial Z_1} \cdot \frac{\partial Z_1}{\partial w_1}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial Z_2} \cdot \frac{\partial Z_2}{\partial A_1} \cdot \frac{\partial A_1}{\partial Z_1} \cdot \frac{\partial Z_1}{\partial b_1}$$

$$\frac{\partial Z_2}{\partial A_1} = W_2^T, \quad \frac{\partial A_1}{\partial Z_1} = A_1(1-A_1)$$

$$L = -(y \log A_2 + (1-y) \log (1-A_2))$$

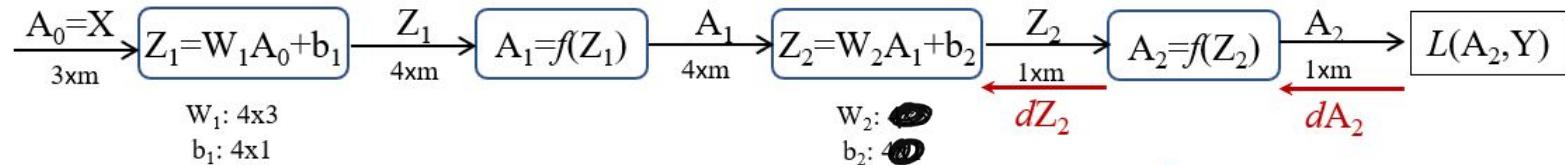
$$\frac{\partial Z_1}{\partial w_1} = A_0^T, \quad \frac{\partial Z_1}{\partial b_1} = [1, \dots, 1]^T$$

$$\frac{\partial L}{\partial A_2} = \frac{A_2 - y}{A_2(1-A_2)}, \quad \frac{\partial Z_2}{\partial w_2} = A_1^T$$

$$\frac{\partial A_2}{\partial Z_2} = A_2(1-A_2), \quad \frac{\partial Z_2}{\partial b_2} = [1, 1, \dots, 1]^T$$

1xm 4x1 4xm 4mx1

# MLP Learning - Back Propagation



$$\frac{\partial L}{\partial w_2} = \frac{(A_2 - Y)}{A_2(1-A_2)} \cdot \underbrace{A_2(1-A_2)}_{1 \times 4} \cdot \underbrace{A_1^T}_{m \times 4}, \quad \frac{\partial L}{\partial b_2} = \frac{(A_2 - Y)}{A_2(1-A_2)} \cdot \underbrace{A_2(1-A_2)}_{1 \times 1} \cdot \underbrace{[1 \dots 1]^T}_{m \times 1}$$

$$\frac{\partial L}{\partial w_1} = \frac{(A_2 - Y)}{A_2(1-A_2)} \cdot \underbrace{A_2(1-A_2)}_{1 \times 1} \cdot \underbrace{(w_2)^T}_{4 \times 1} \cdot \underbrace{A_1}_{1 \times 4} \cdot \underbrace{A_0^T}_{m \times 3}$$

$$\frac{\partial L}{\partial w_1} = \underbrace{(w_2)^T}_{4 \times 1} \underbrace{(A_2 - Y)}_{1 \times m} \underbrace{A_1}_{4 \times m} \underbrace{(1-A_1)}_{4 \times m} \underbrace{A_0^T}_{m \times 3}$$

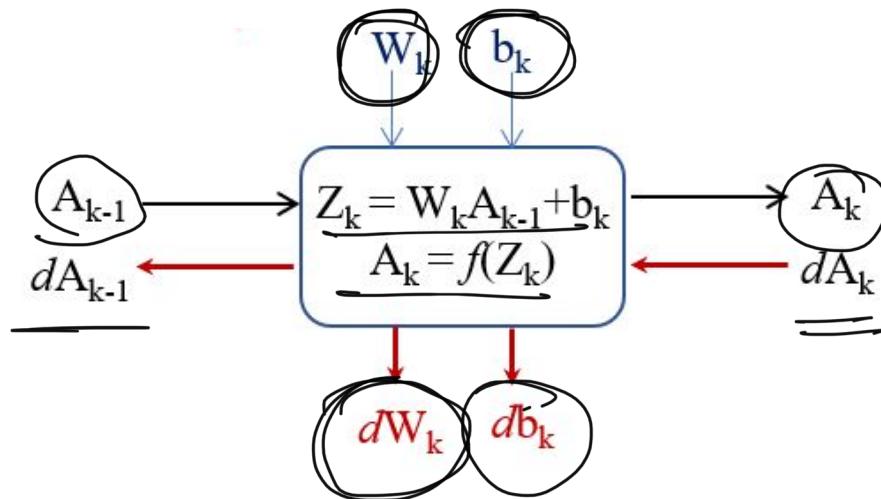
elementwise multiplication

$$\frac{\partial L}{\partial b_1} = \frac{(A_2 - Y)}{A_2(1-A_2)} \cdot \underbrace{A_2(1-A_2)}_{1 \times 1} \cdot \underbrace{w_2^T}_{4 \times 1} \underbrace{A_1}_{1 \times 4} \underbrace{(1-A_1)}_{1 \times m} \underbrace{[1 \dots 1]^T}_{m \times 1}$$

$$\frac{\partial L}{\partial b_1} = \underbrace{w_2^T}_{4 \times 1} \underbrace{(A_2 - Y)}_{1 \times m} \underbrace{A_1}_{4 \times m} \underbrace{(1-A_1)}_{4 \times m} \underbrace{[1 \dots 1]^T}_{m \times 1}$$

elementwise multiplication

# MLP Learning - Back Propagation



# Questions?