

# **COMP47590**

## **ADVANCED MACHINE LEARNING**

### **DEEP LEARNING - ANNs 2**

Dr. Brian Mac Namee



## **Information**

Email:

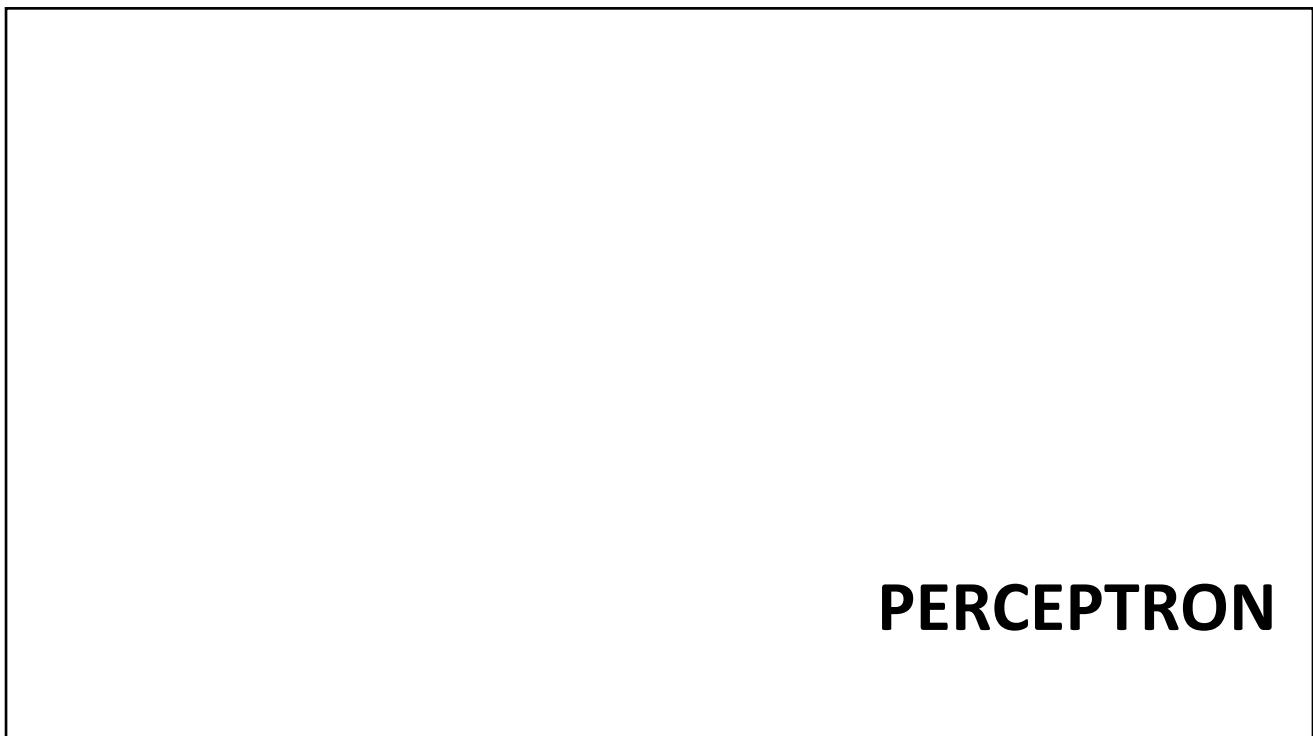
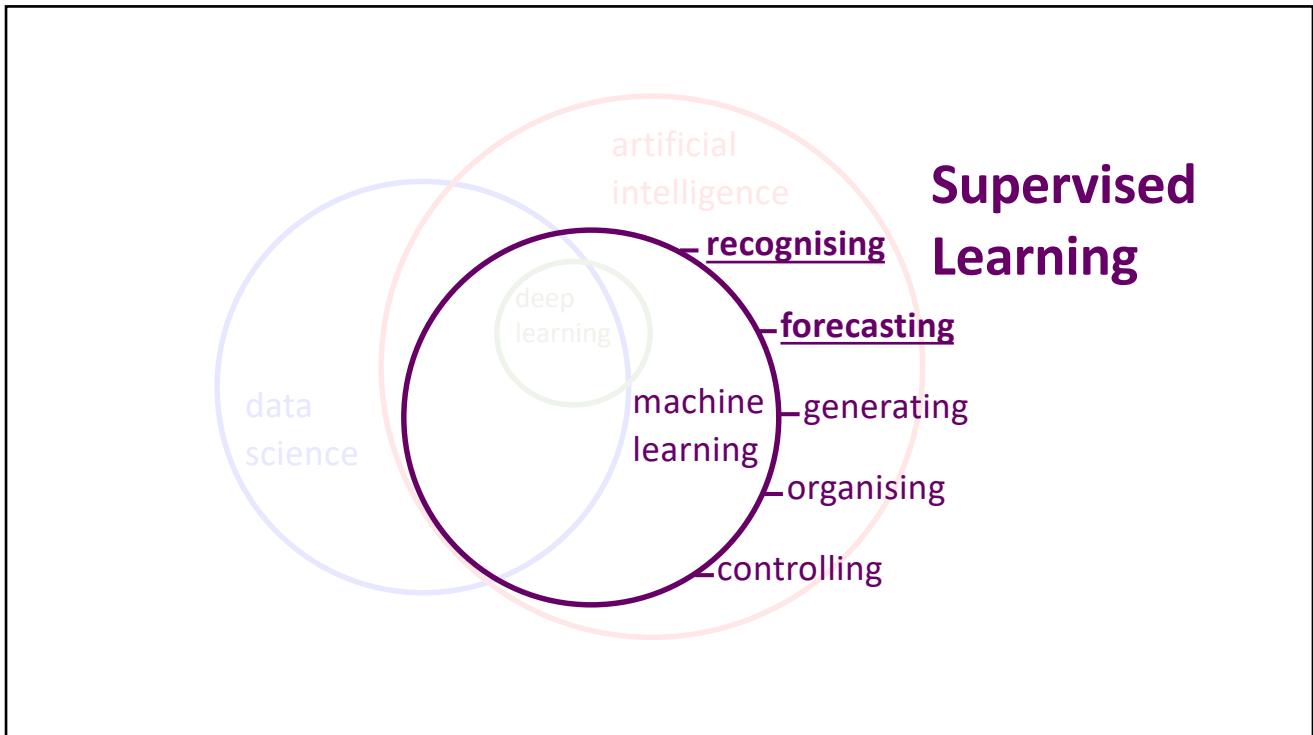
[Brian.MacNamee@ucd.ie](mailto:Brian.MacNamee@ucd.ie)

Course Materials:

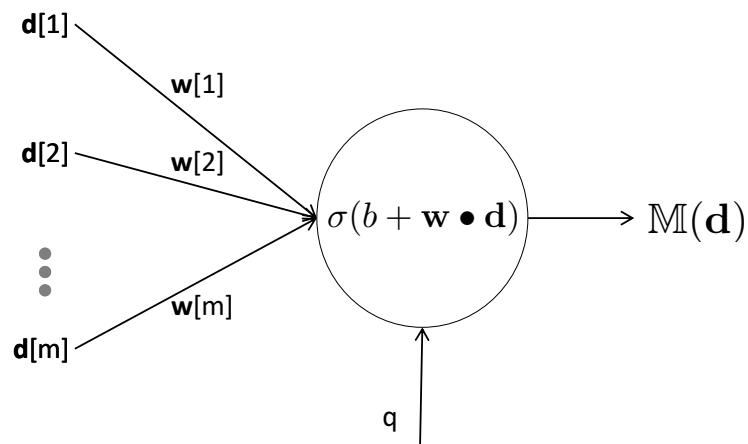
All material posted on UCD CS moodle

<https://csmoodle.ucd.ie/moodle/course/view.php?id=756>

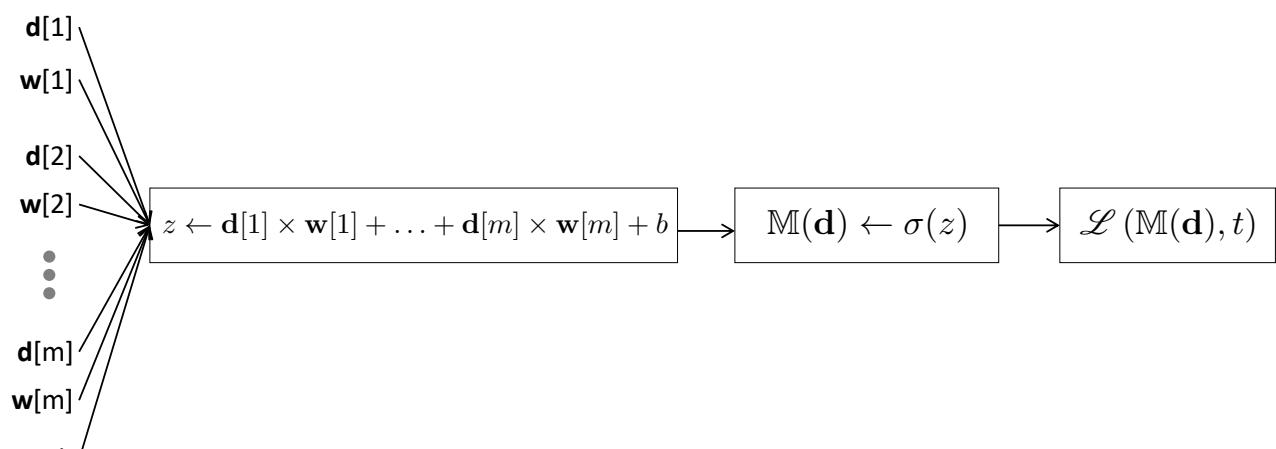
Enrolment key **UCDAvML2019**



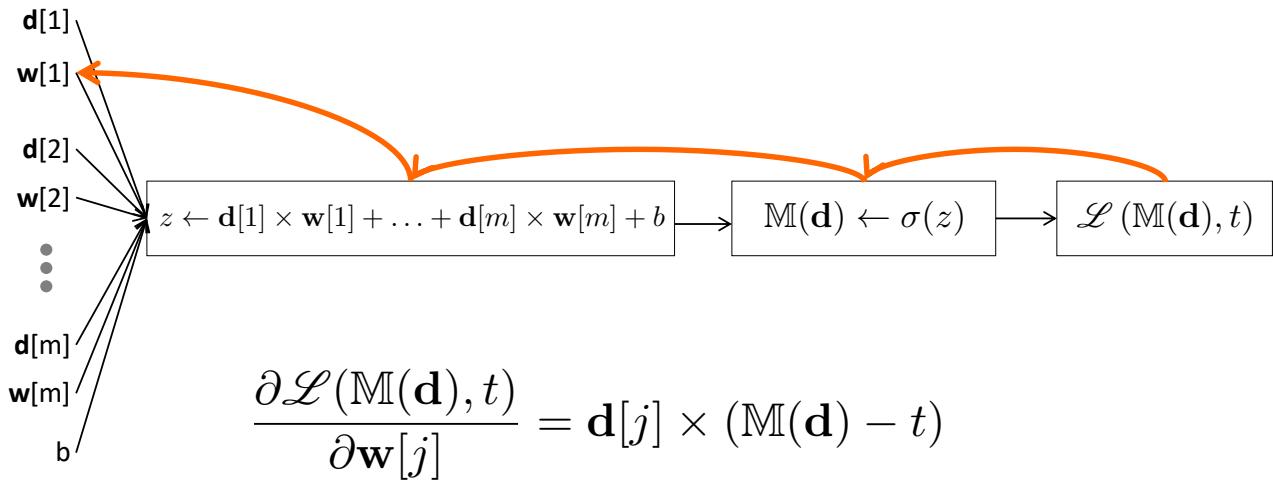
## Perceptron



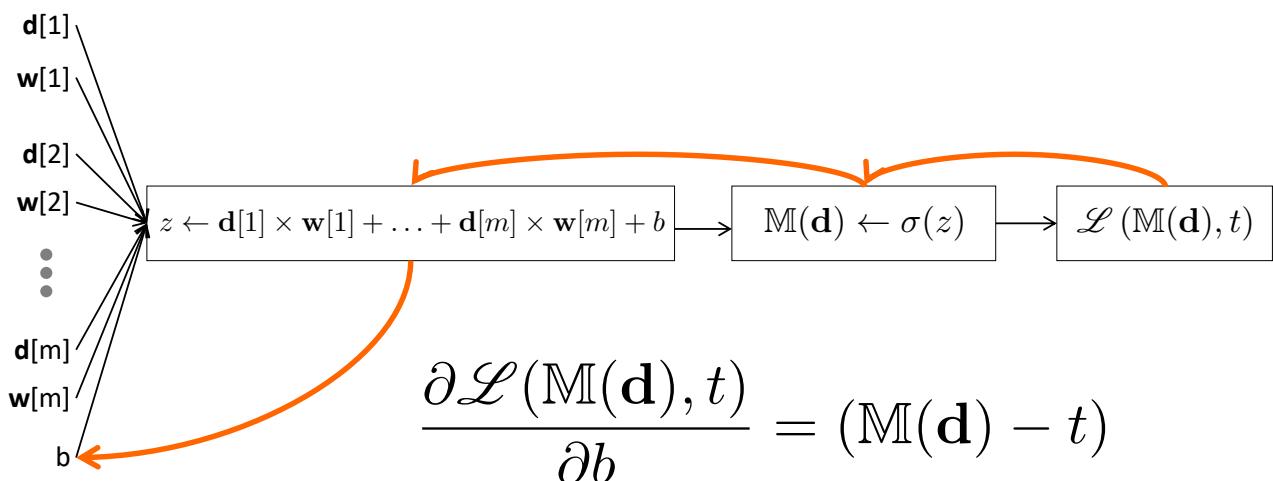
## Perceptron



## Perceptron

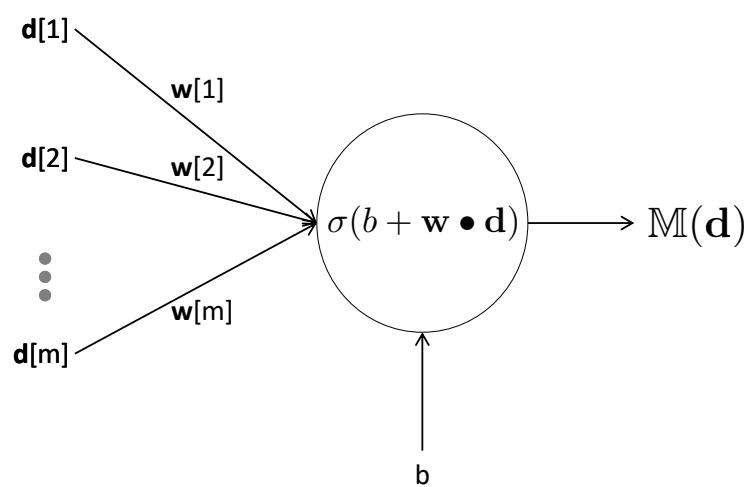


## Perceptron

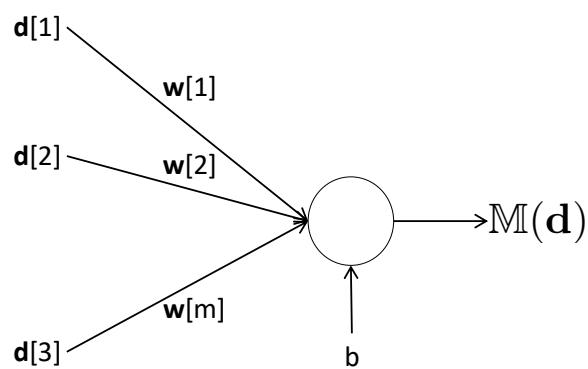


# ARTIFICIAL NEURAL NETWORKS

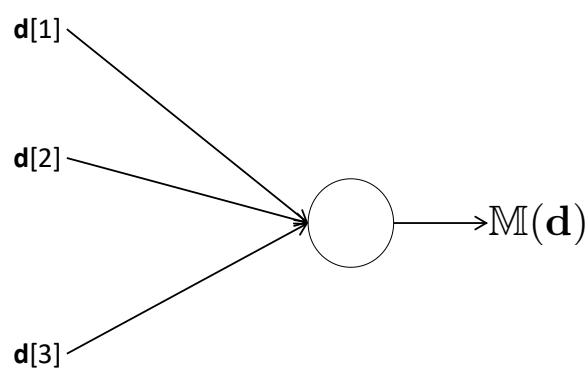
## Perceptron



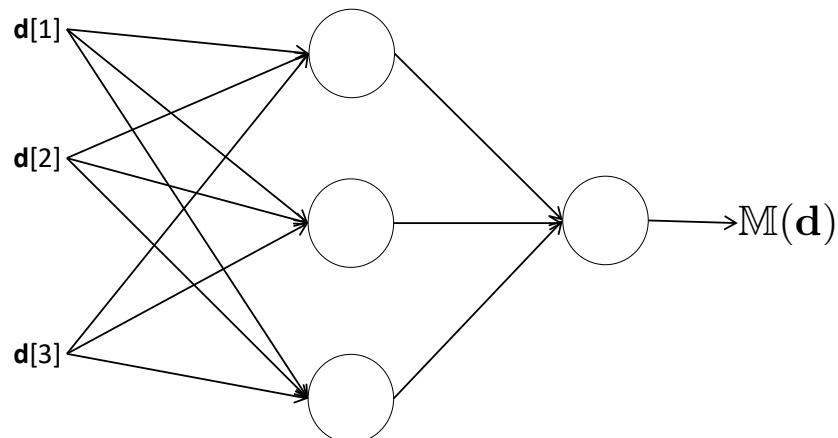
## Perceptron



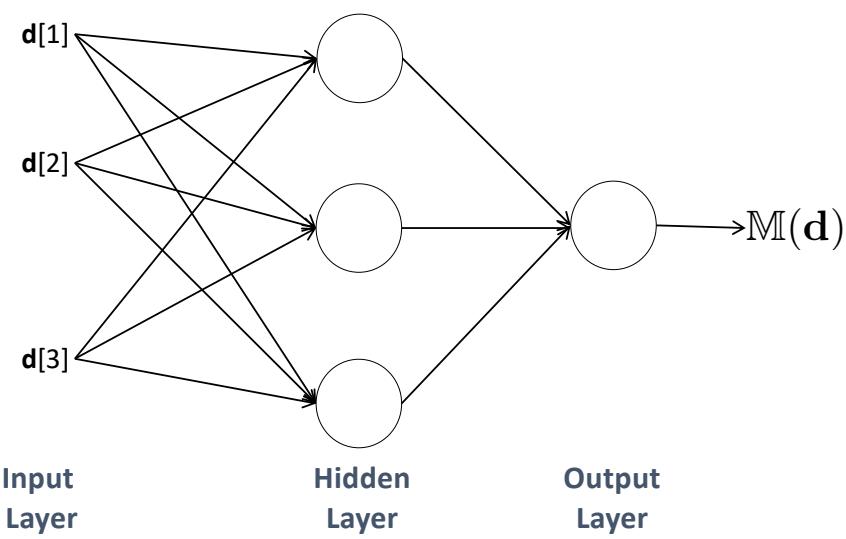
## Perceptron



## Multi Layer Perceptron

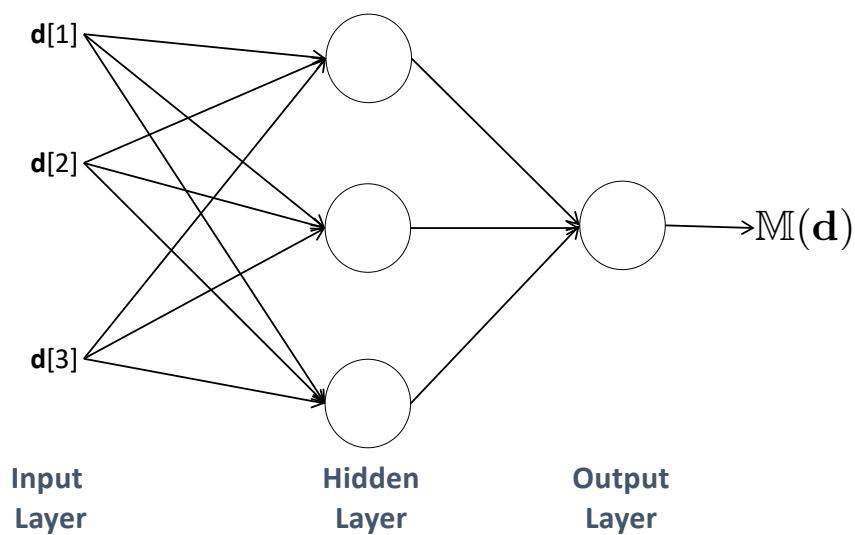


## Multi Layer Perceptron



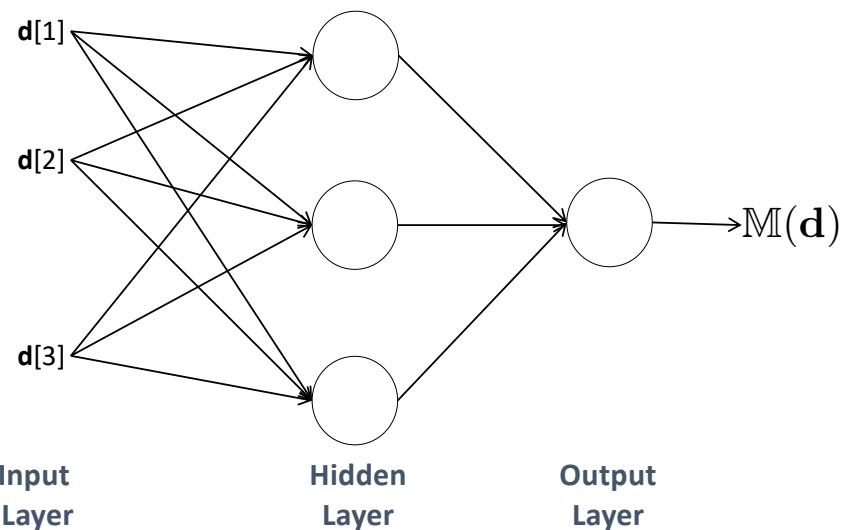
## BACKPROPOGATION OF ERRORS

### Multi Layer Perceptron



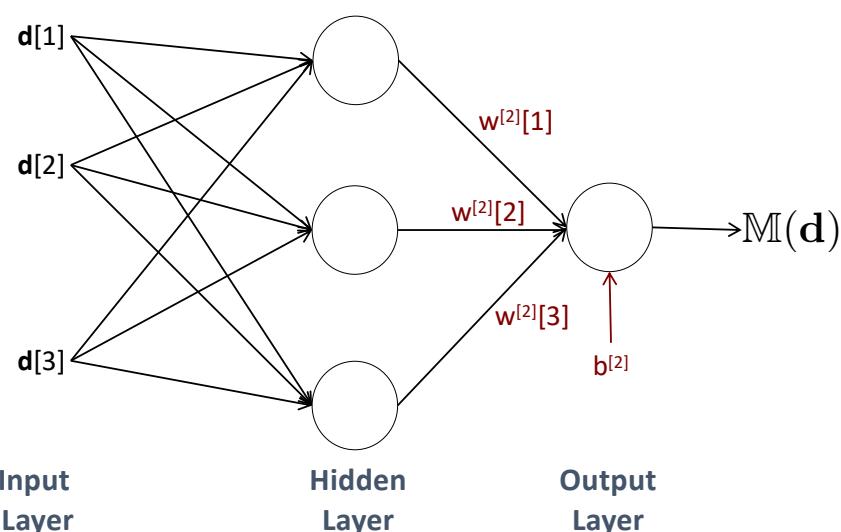
## Multi Layer Perceptron

If we remember what is going on at the output layer



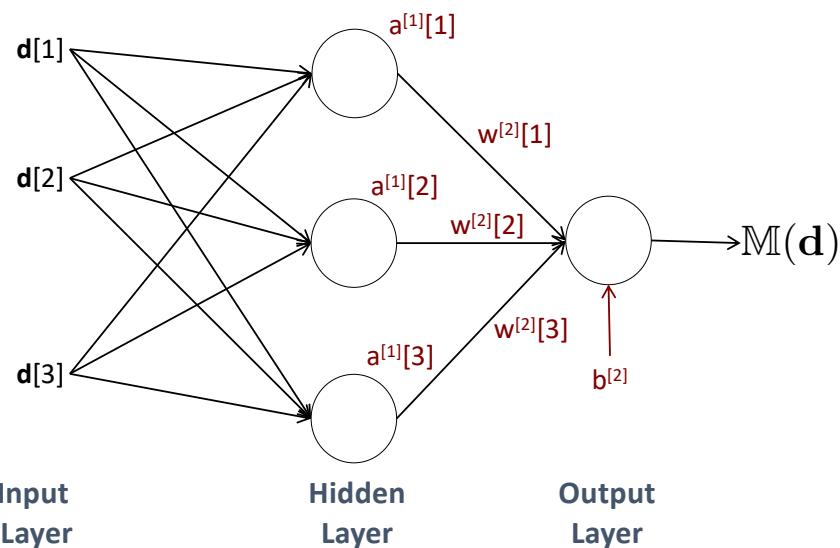
## Multi Layer Perceptron

If we remember what is going on at the output layer



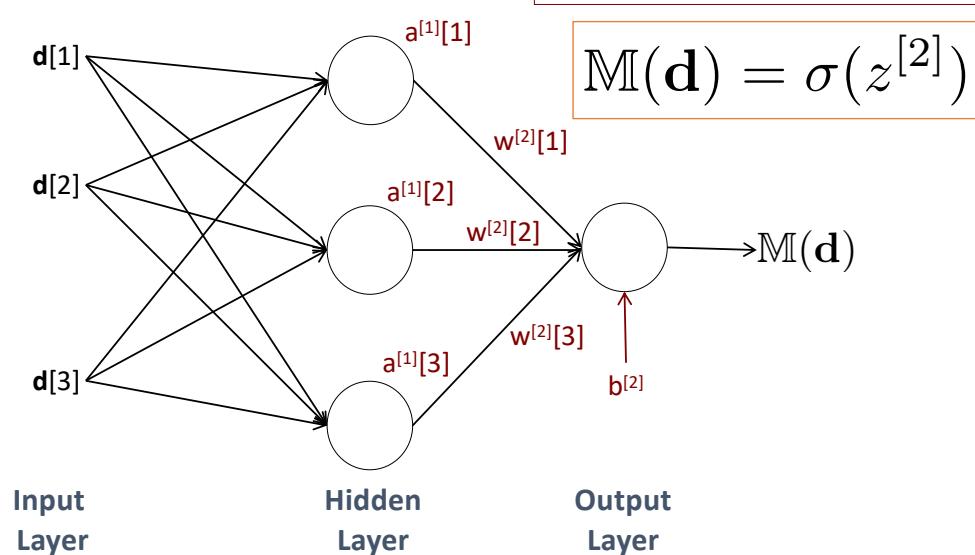
## Multi Layer Perceptron

If we remember what is going on at the output layer



## Multi Layer Perceptron

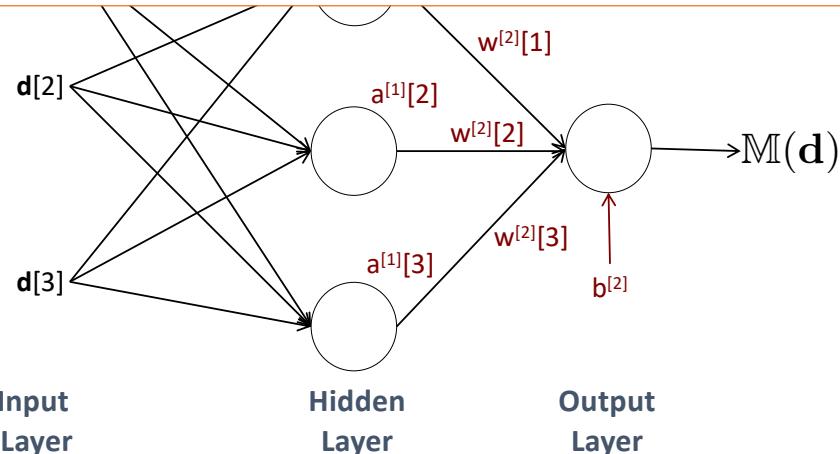
If we remember what is going on at the output layer



## Multi Layer Perceptron

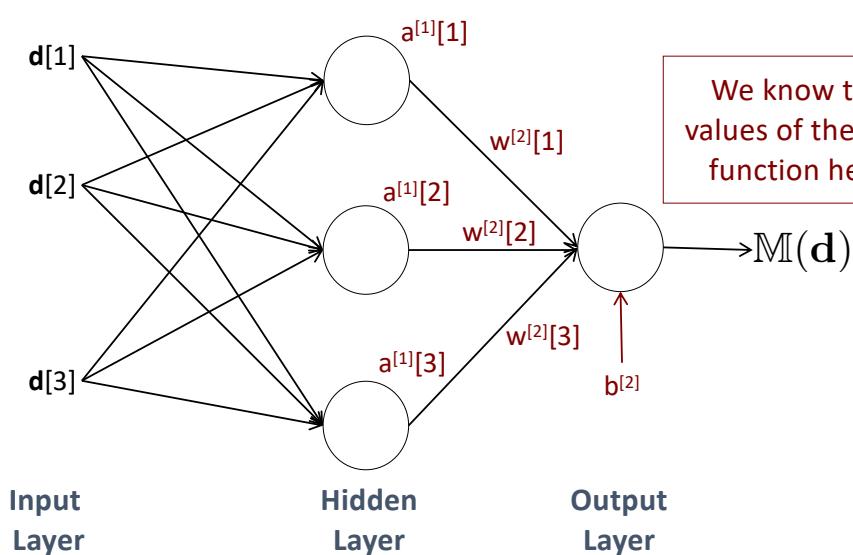
If we remember what is going on at the output layer

$$z^{[2]} = \mathbf{w}^{[2]}[1] \times \mathbf{a}^{[1]}[1] + \mathbf{w}^{[2]}[2] \times \mathbf{a}^{[1]}[2] + \mathbf{w}^{[2]}[3] \times \mathbf{a}^{[1]}[3]$$

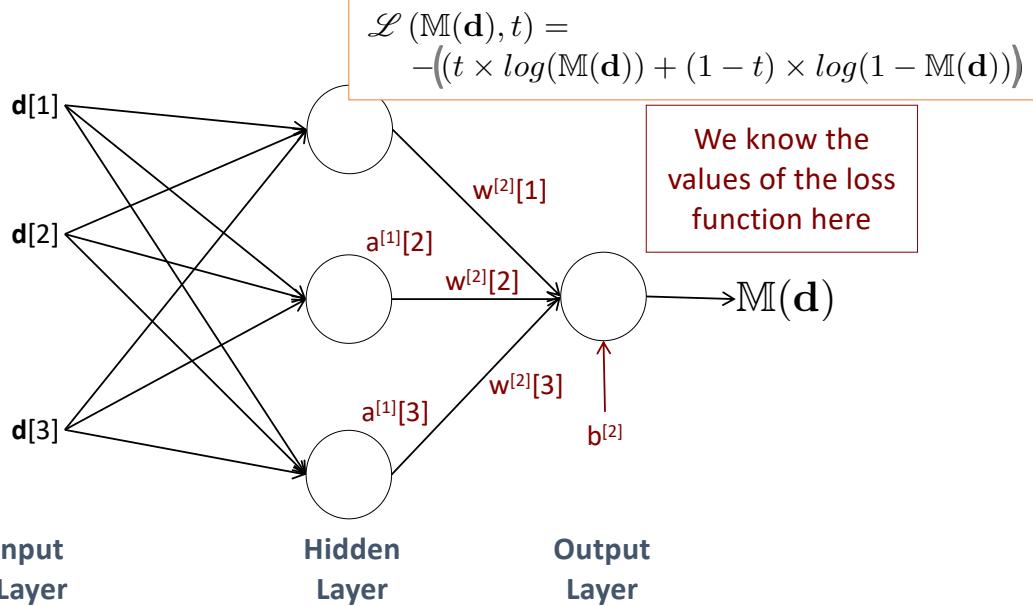


## Multi Layer Perceptron

We know the values of the loss function here

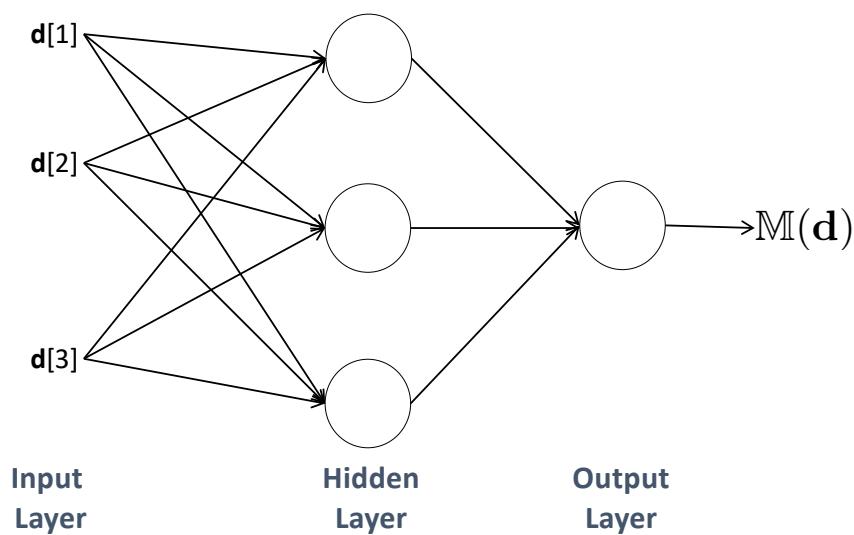


## Multi Layer Perceptron



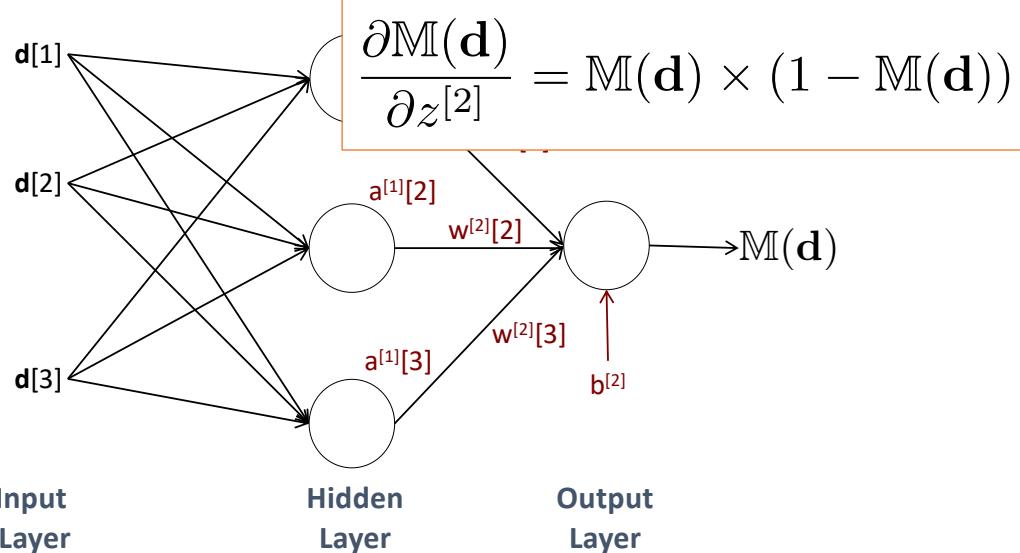
## Multi Layer Perceptron

So we can calculate the derivative of the loss function with respect to the weights coming into this unit



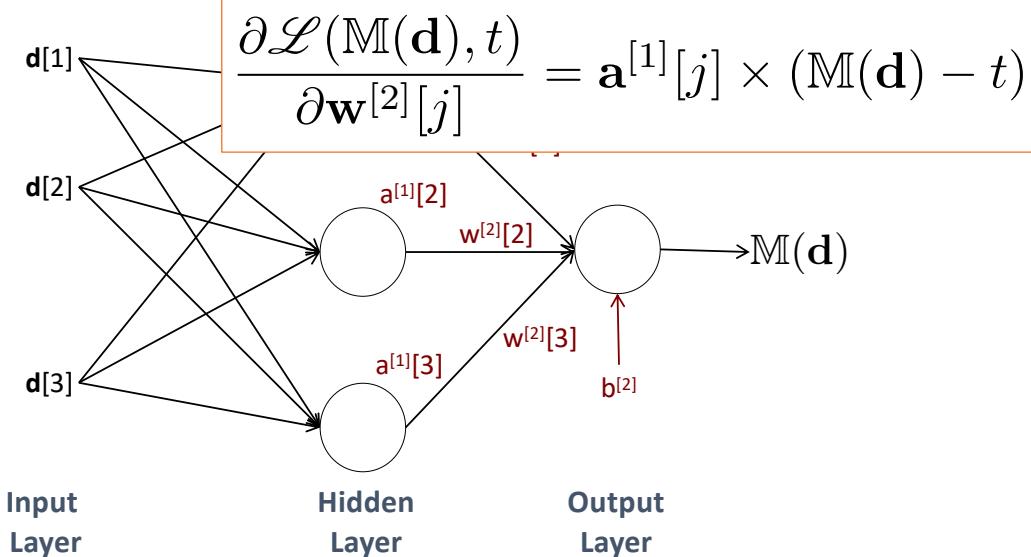
## Multi Layer Perceptron

So we can calculate the derivative of the loss function with respect to the output of the last unit

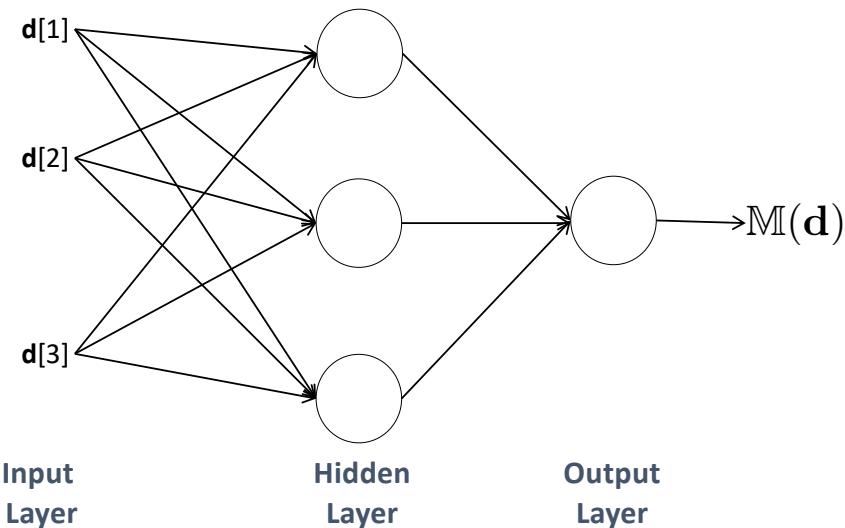


## Multi Layer Perceptron

So we can calculate the derivative of the loss function with respect to the weights coming into this unit

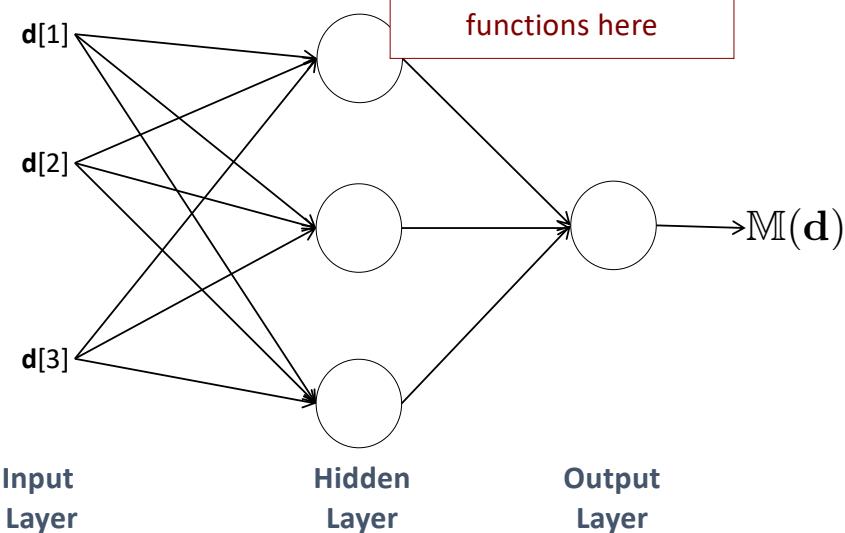


## Multi Layer Perceptron

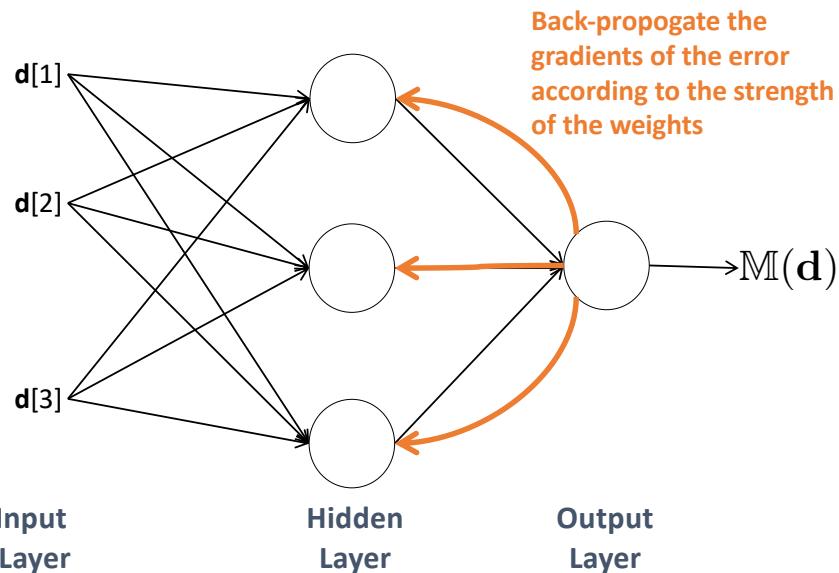


## Multi Layer Perceptron

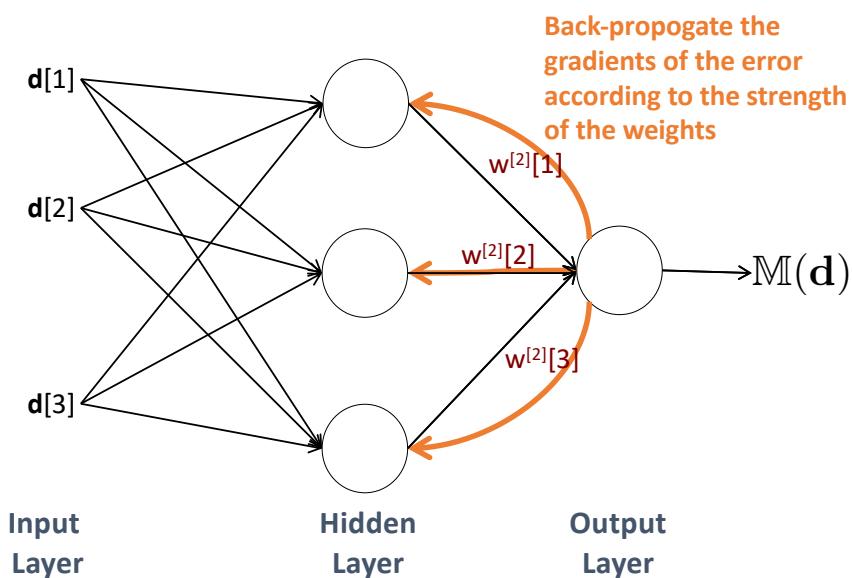
We don't know the value of the cost functions here



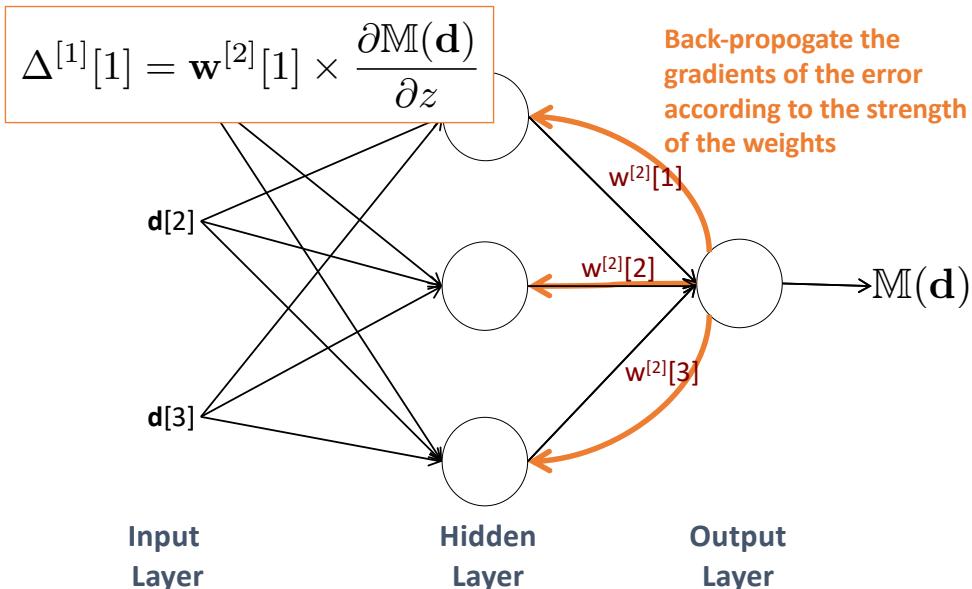
## Multi Layer Perceptron



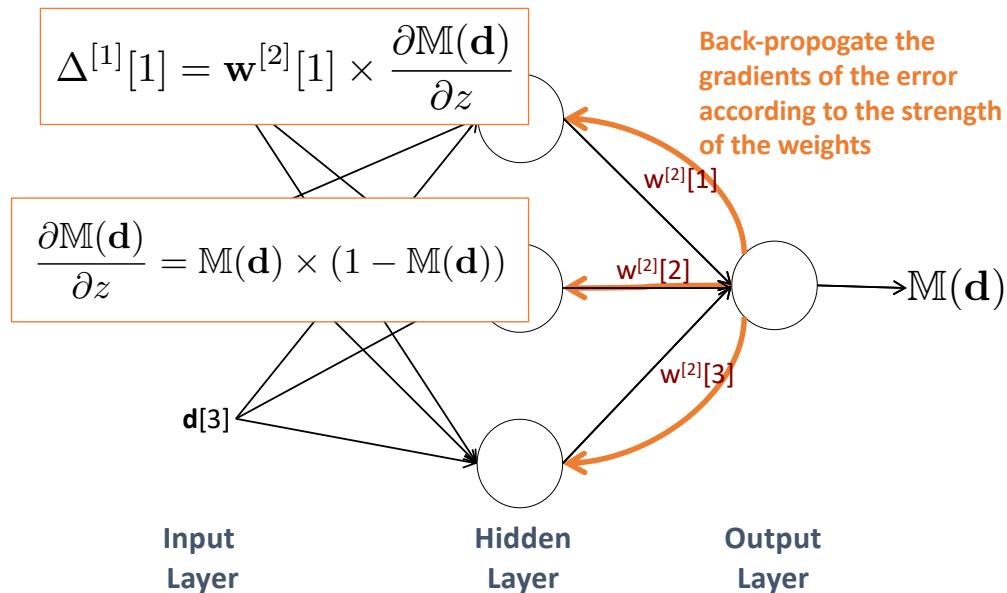
## Multi Layer Perceptron



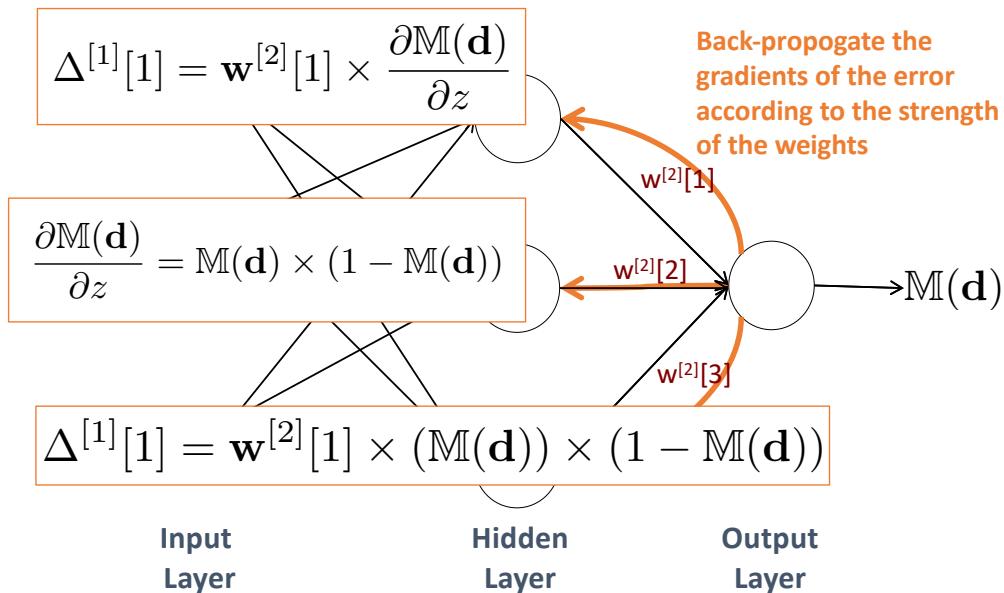
## Multi Layer Perceptron



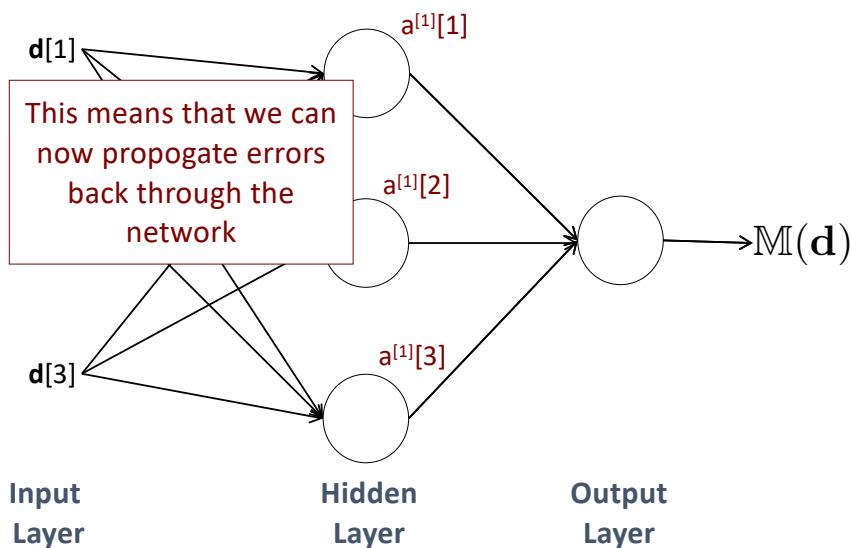
## Multi Layer Perceptron



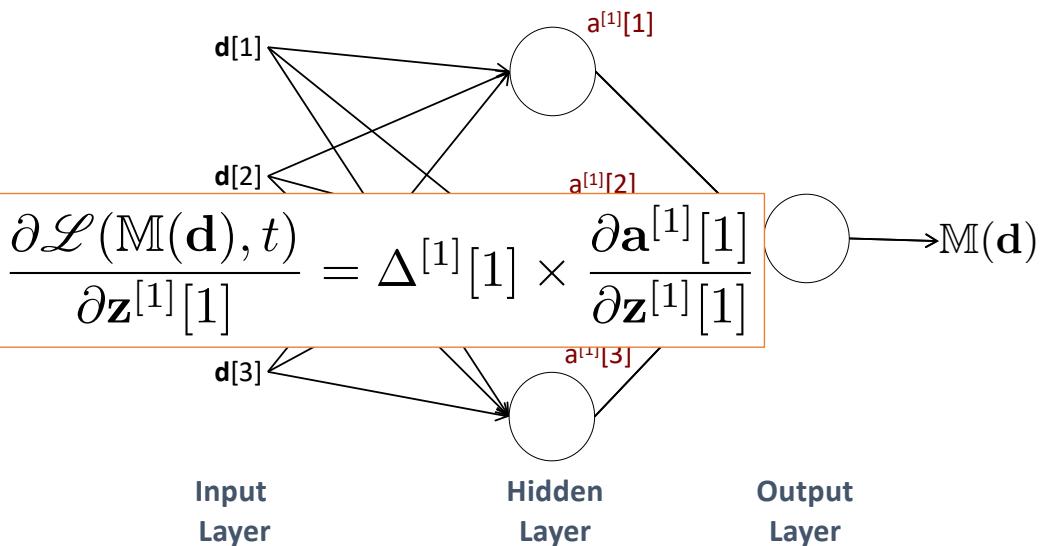
## Multi Layer Perceptron



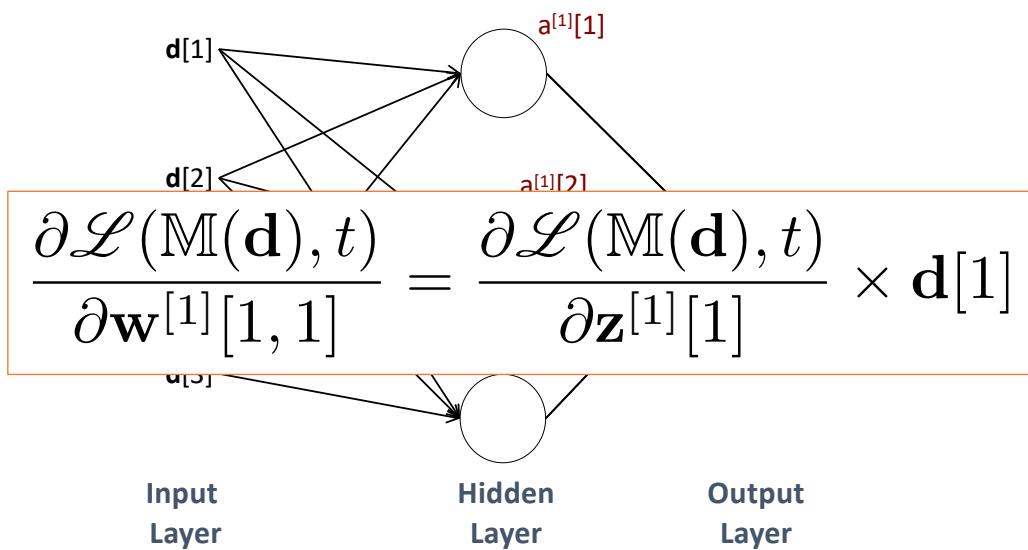
## Multi Layer Perceptron



## Multi Layer Perceptron

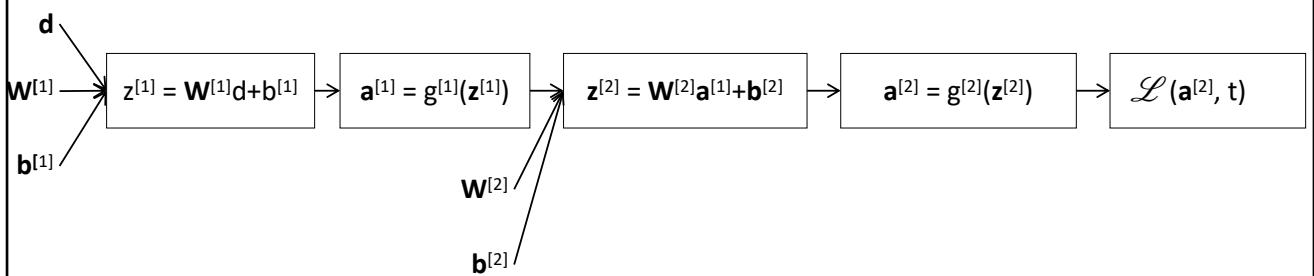


## Multi Layer Perceptron



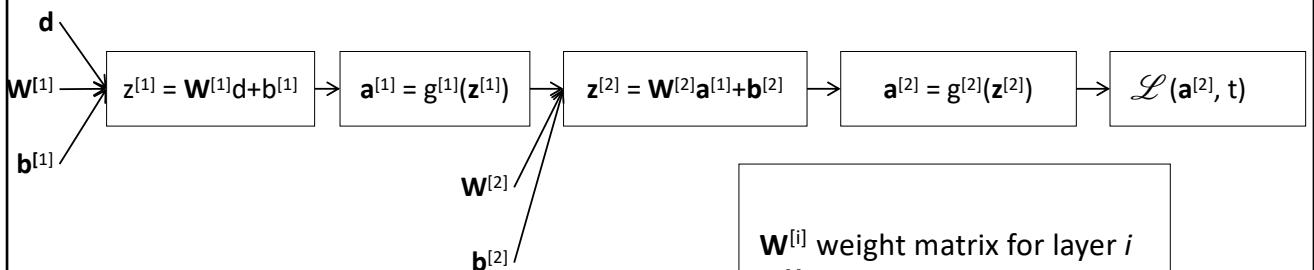
## A NEW NOTATION

### Multi Layer Perceptron



## Multi Layer Perceptron

We adopt a notation here (following Andrew Ng) that makes the layer structure of networks very explicit



$\mathbf{W}^{[i]}$  weight matrix for layer  $i$   
 $\mathbf{b}^{[i]}$  bias vector for layer  $i$

## A New Notation

$$d\mathbf{z}^{[2]} = \mathbf{a}^{[2]} - \mathbf{t}$$

$$d\mathbf{W}^{[2]} = d\mathbf{z}^{[2]}\mathbf{a}^{[1]\top}$$

$$d\mathbf{b}^{[2]} = d\mathbf{z}^{[2]}$$

$$d\mathbf{z}^{[1]} = \mathbf{W}^{[2]\top}d\mathbf{z}^{[2]} * g^{[1]'}(\mathbf{z}^{[1]})$$

$$d\mathbf{W}^{[1]} = d\mathbf{z}^{[1]}\mathbf{d}^\top$$

$$d\mathbf{b}^{[1]} = d\mathbf{z}^{[1]}$$

where

$$d\mathbf{X} = \frac{\partial J(\mathbf{w}, b)}{\partial \mathbf{X}}$$

## GRADIENT DESCENT

### Gradient Descent

We only need to update our weight update rules to take account of the extra layers

$$\mathbf{W}^{[i]} = \mathbf{W}^{[i]} - \alpha \mathbf{dW}^{[i]}$$

$$\mathbf{b}^{[i]} = \mathbf{b}^{[i]} - \alpha \mathbf{db}^{[i]}$$

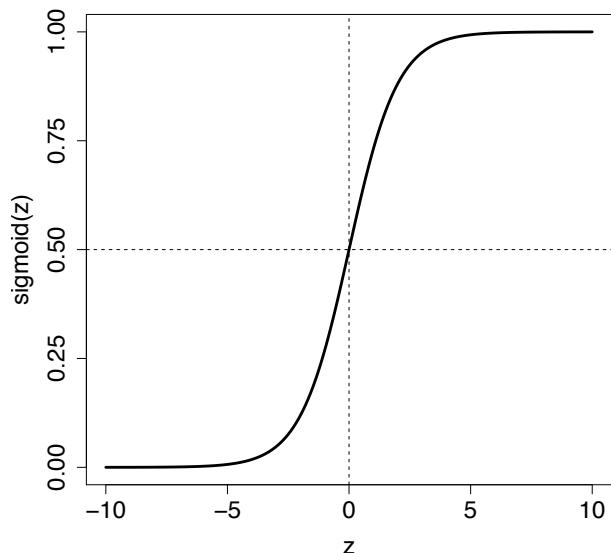
# ACTIVATION FUNCTIONS

## Activation Functions

Unless we introduce a non-linear activation function our networks, no matter how deep, will only be able to learn linear functions

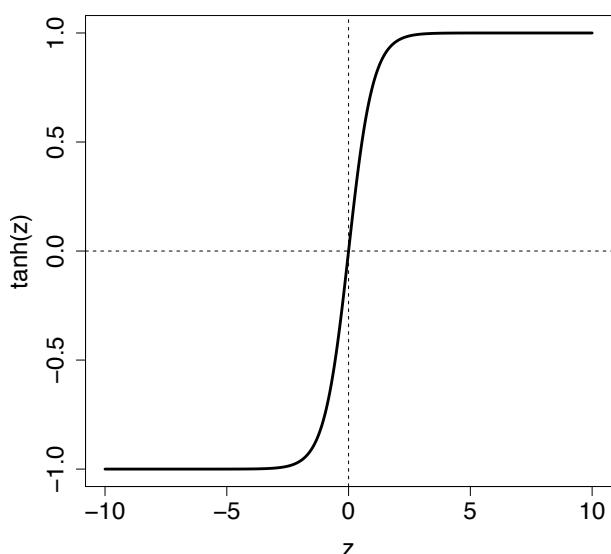
**Sigmoid** activation units used to be most common  
Now replaced with modern alternatives **tanh**,  
**rectified linear units (relu)**, and **leaky relu**

## Activation Functions



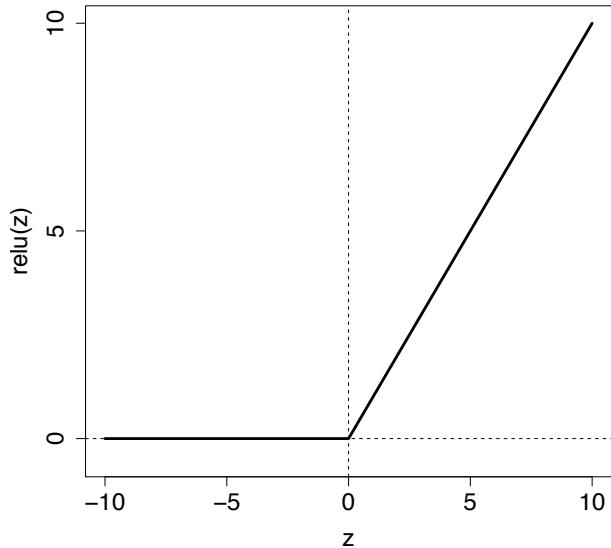
$$\text{sigmoid}(z) = \frac{1}{1 + e^{-z}}$$

## Activation Functions



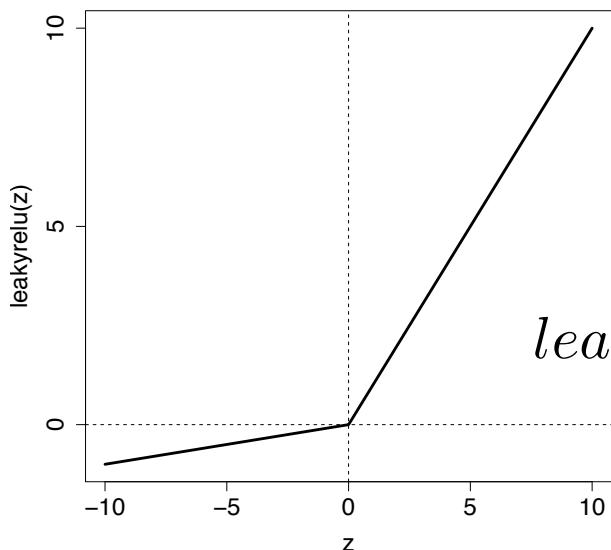
$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

## Activation Functions



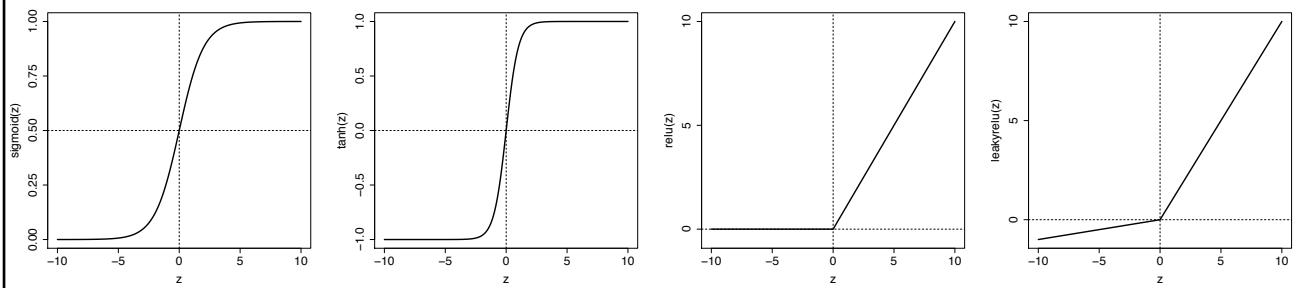
$$\text{relu}(z) = \max(0, z)$$

## Activation Functions



$$\text{leakyrelu}(z) = \max(0.1 z, z)$$

## Activation Functions



**SUMMARY**

## Summary

Expanding the perceptron approach to multiple layers of computational units is straight-forward

The gradient descent algorithm does not require any major changes

The key is the **backpropagation of error** algorithm which allows us to propagate the gradient of the error from the output layer back through the earlier layers in a network

## Questions

