



COMP47590

ADVANCED MACHINE LEARNING

DEEP LEARNING - ANNs 3

Dr. Brian Mac Namee



Information

Email:

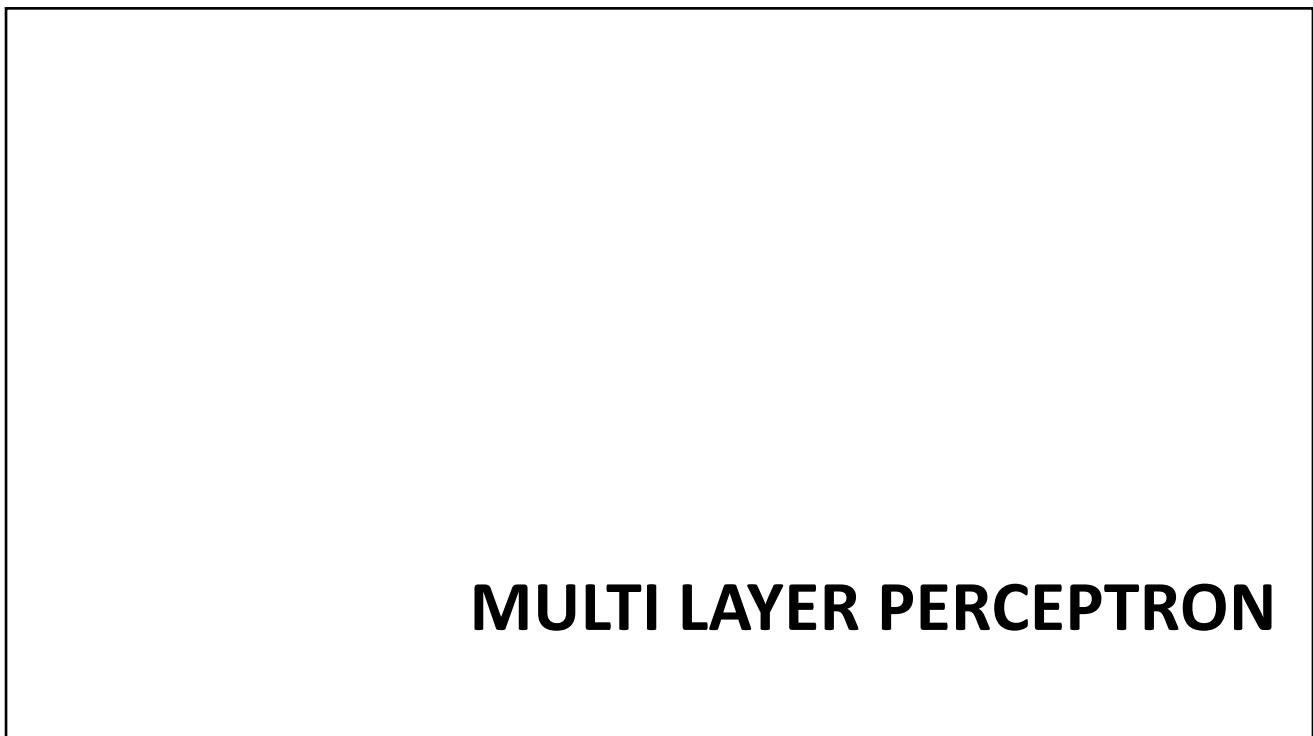
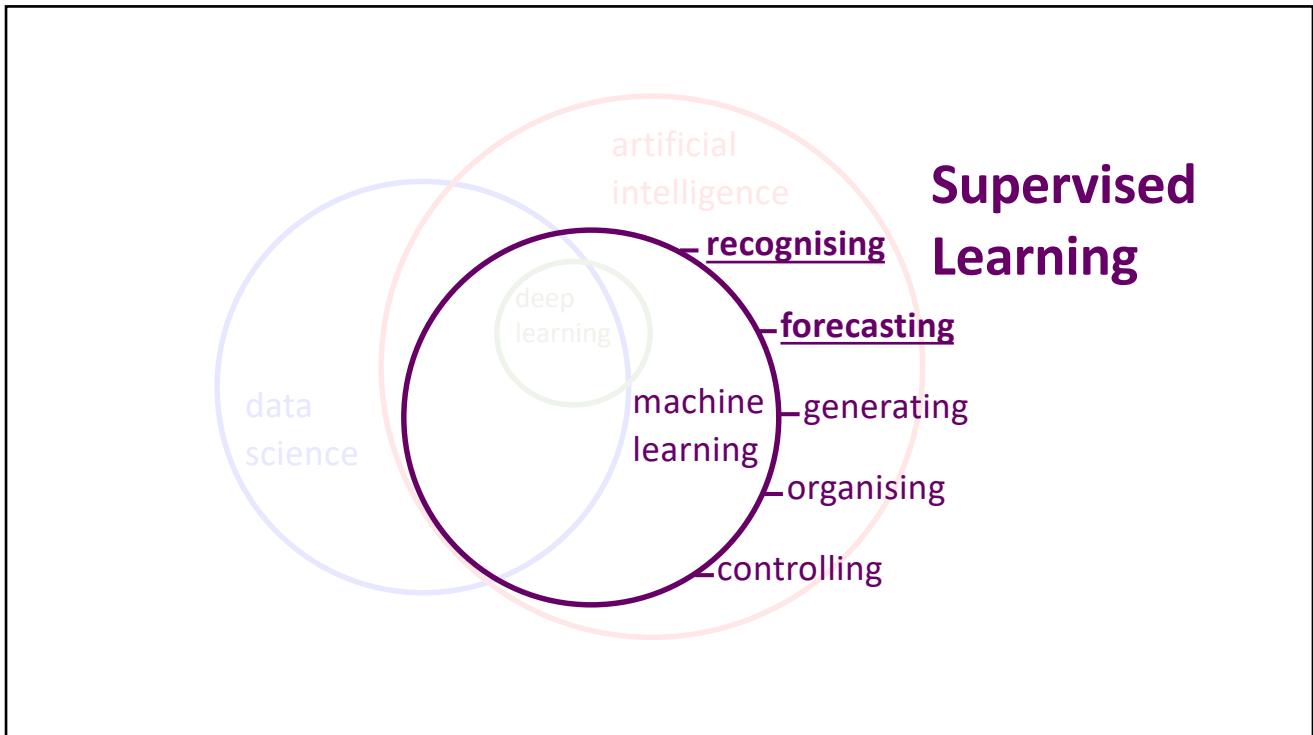
Brian.MacNamee@ucd.ie

Course Materials:

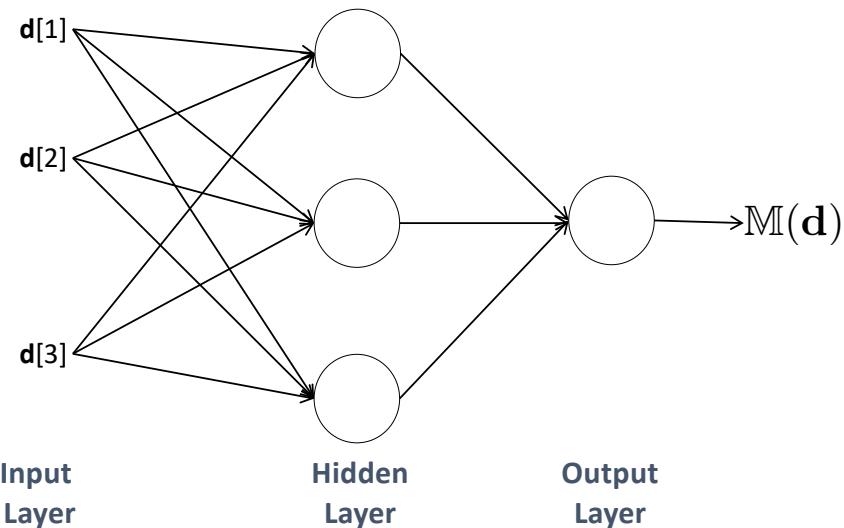
All material posted on UCD CS moodle

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

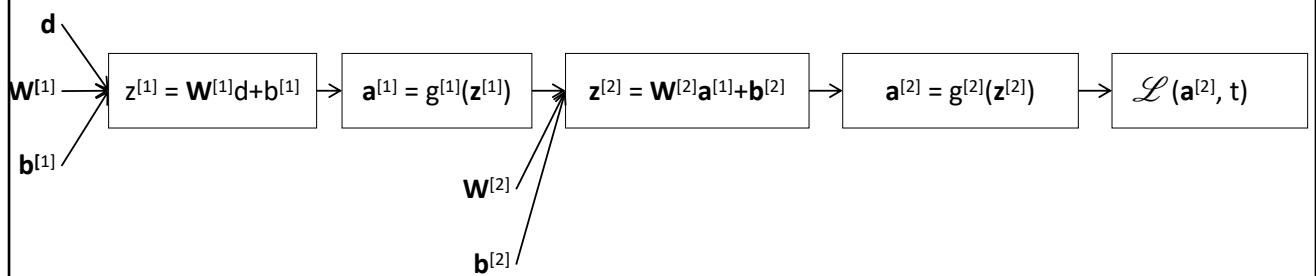
Enrolment key **UCDAvML2019**



Multi Layer Perceptron



MLP: Forward Propagation



MLP: Backward Propagation

$$dz^{[2]} = a^{[2]} - t$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} d^T$$

$$db^{[1]} = dz^{[1]}$$

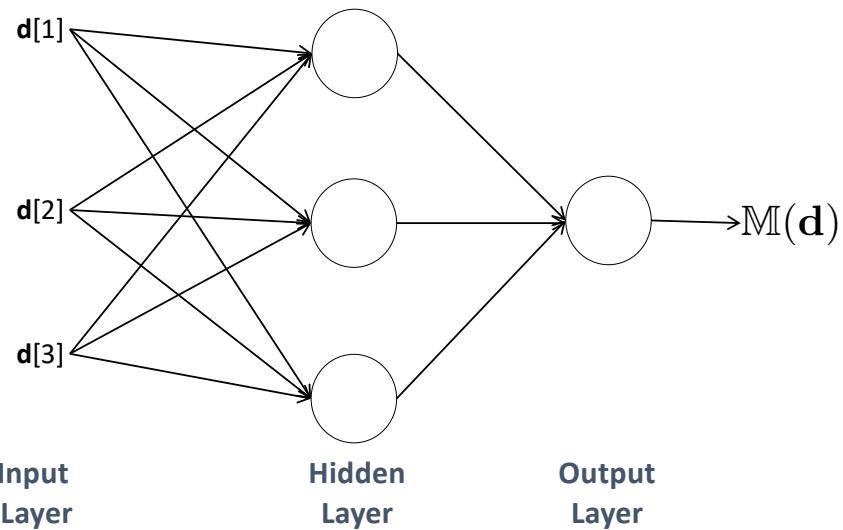
MLP: Gradient Descent

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

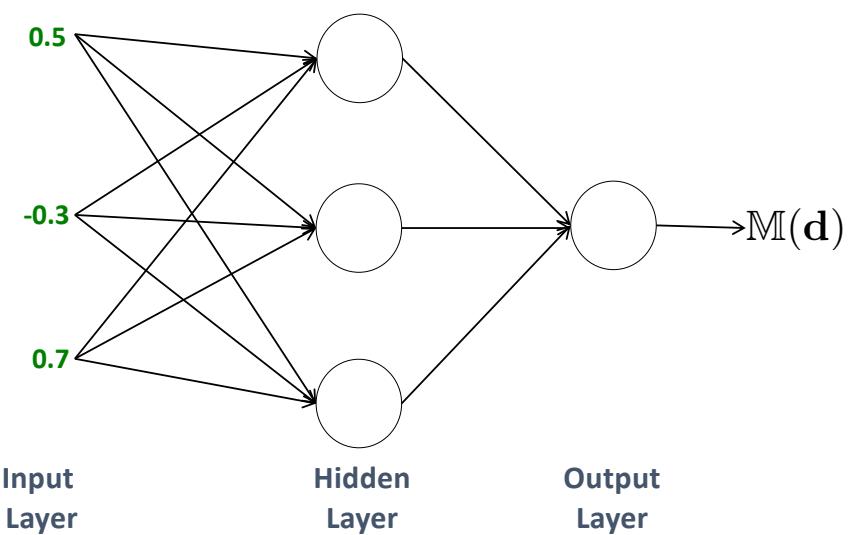
$$W^{[i]} = W^{[i]} - \alpha dW^{[i]}$$

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

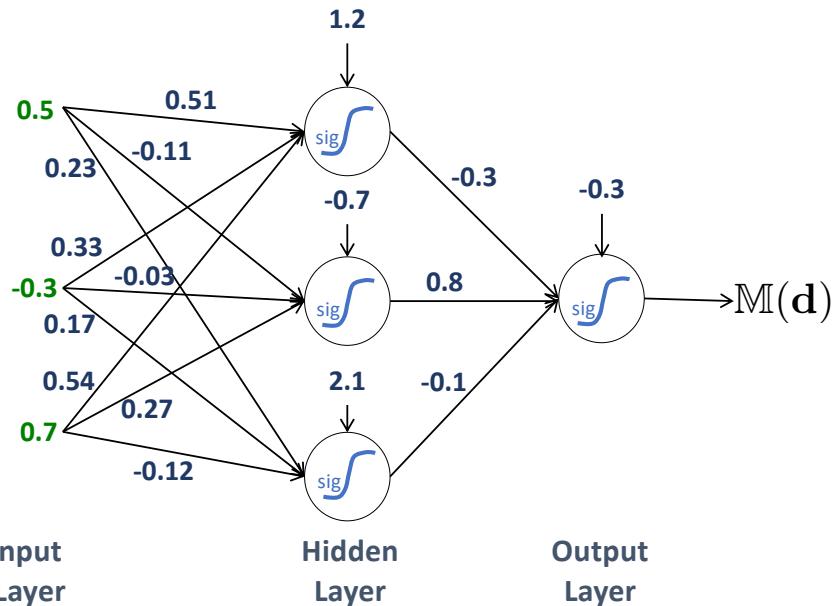
MLP: Exercise



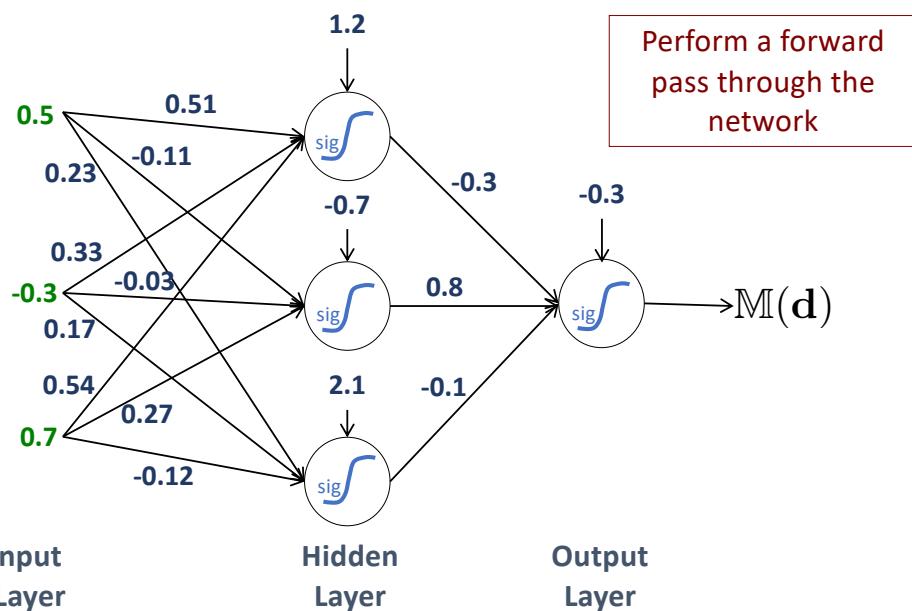
MLP: Exercise



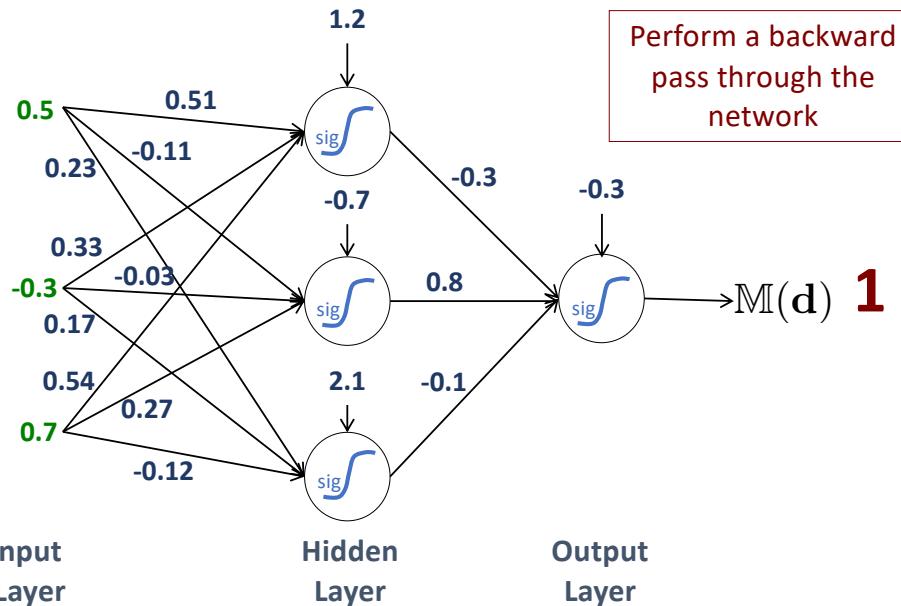
MLP: Exercise



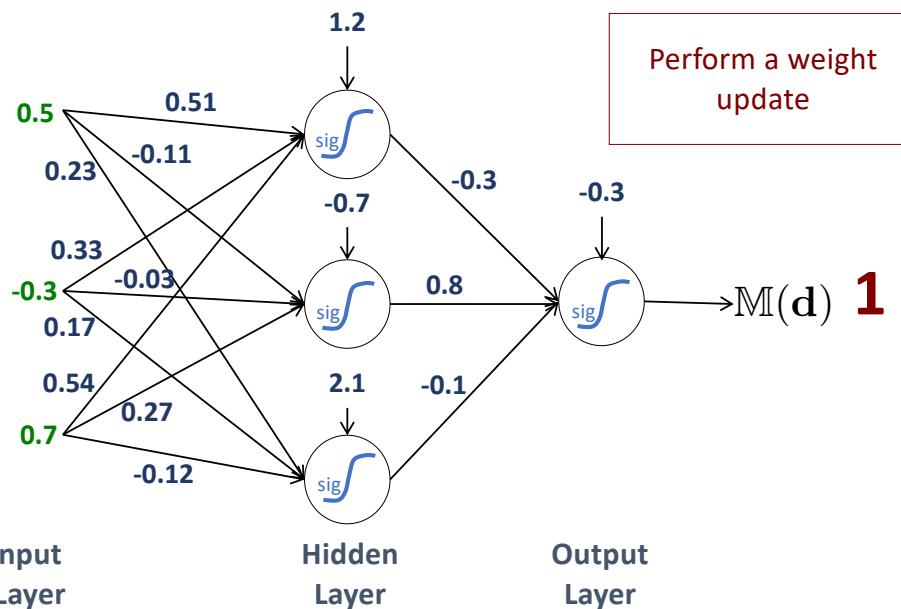
MLP: Exercise



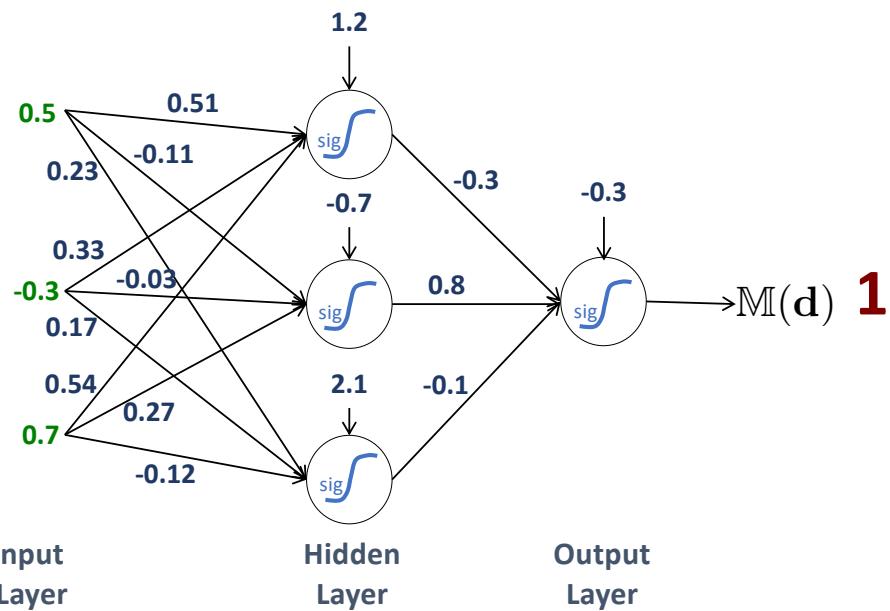
MLP: Exercise



MLP: Exercise

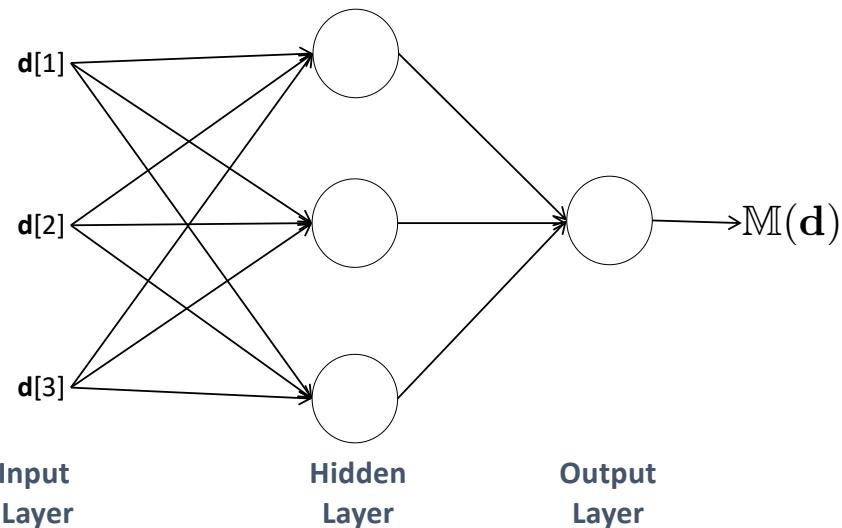


MLP: Exercise

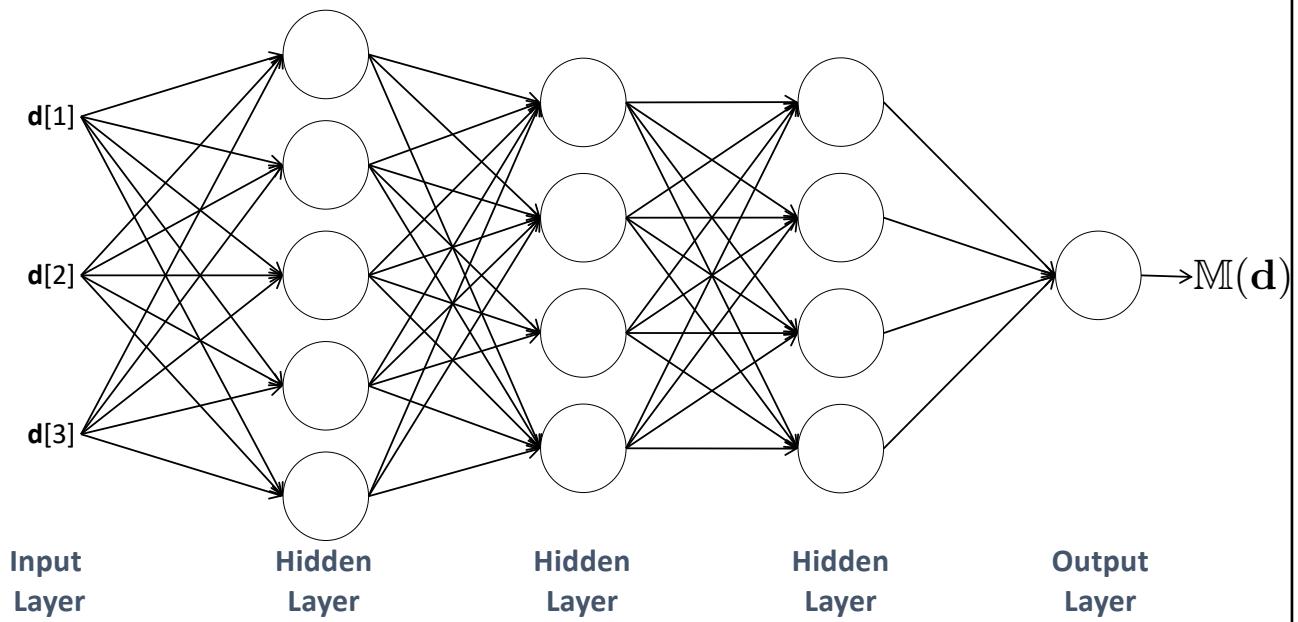


BACKPROPOGATION OF ERRORS

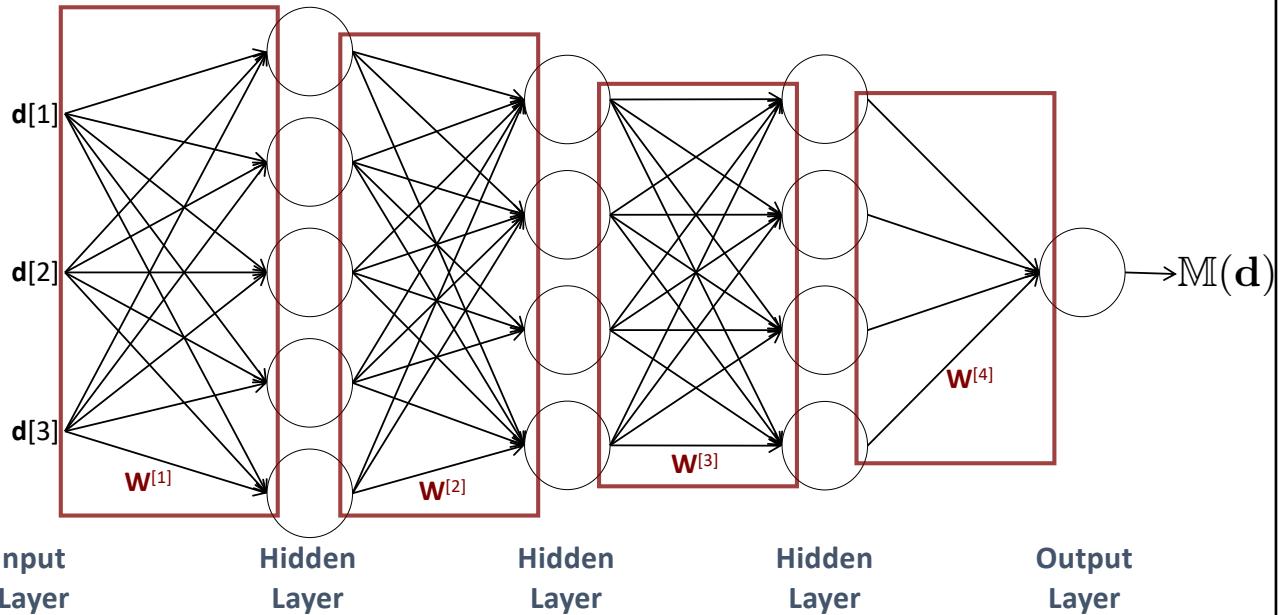
Multi Layer Perceptron



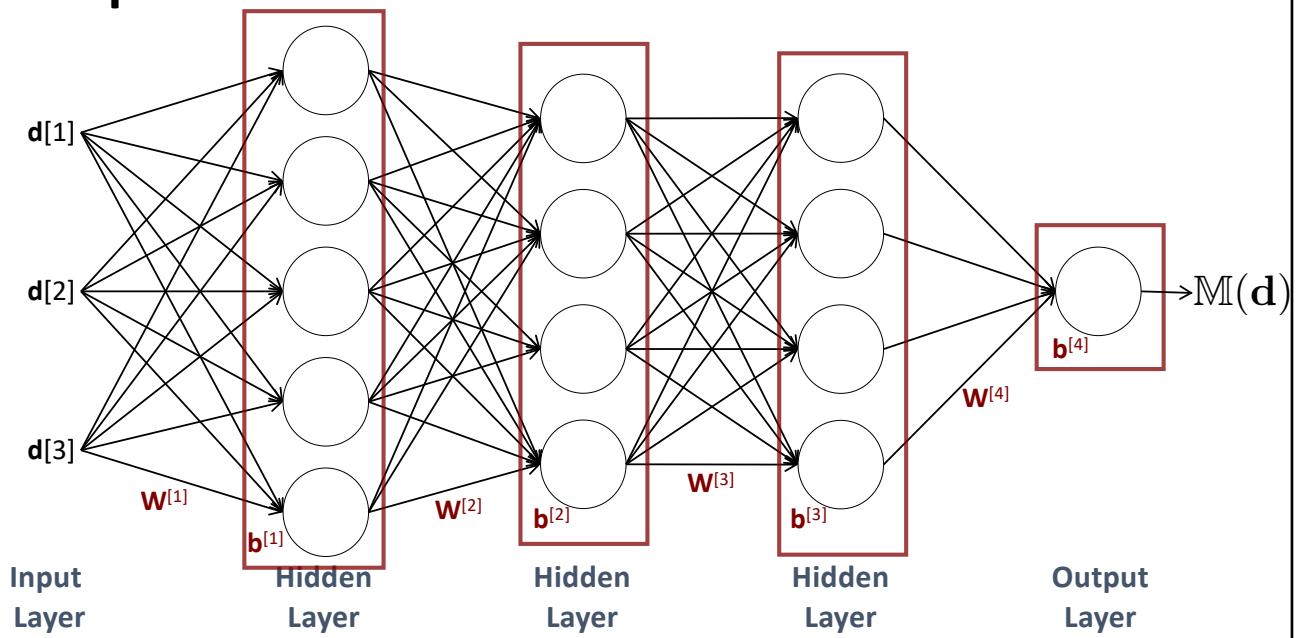
Deep Feed Forward Network

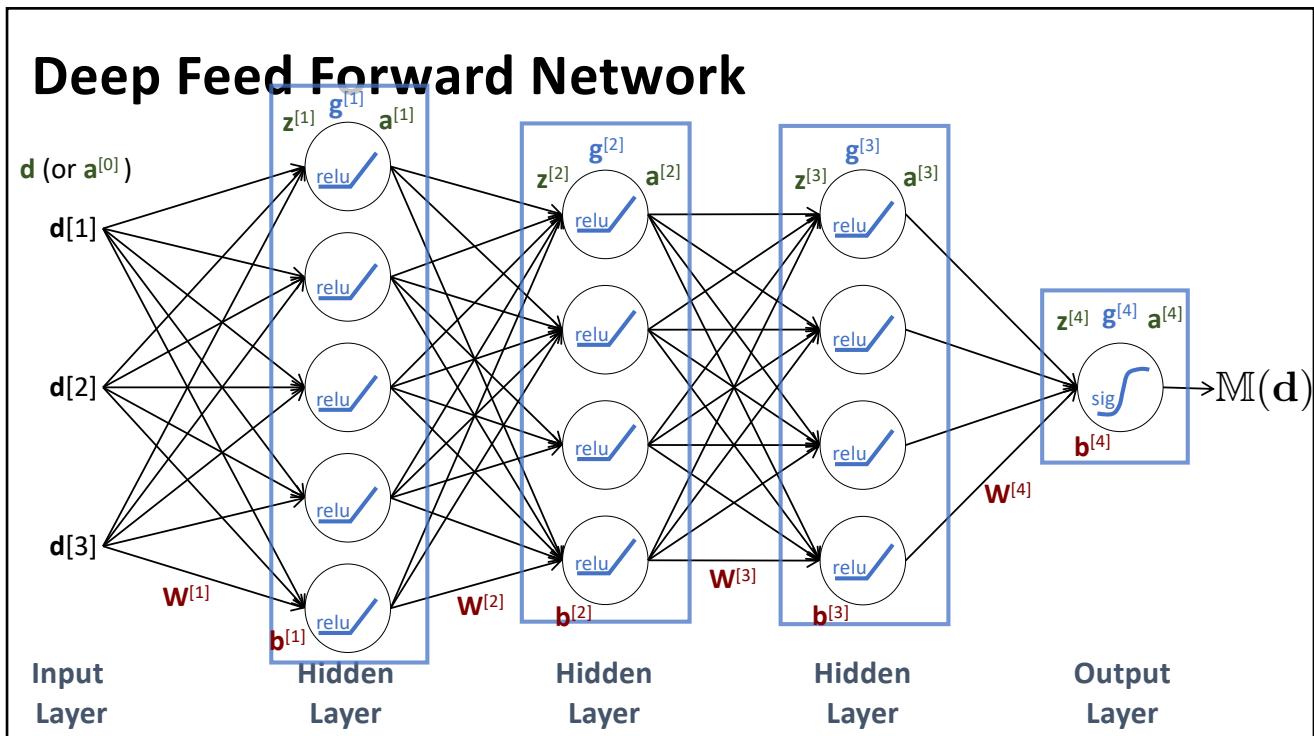
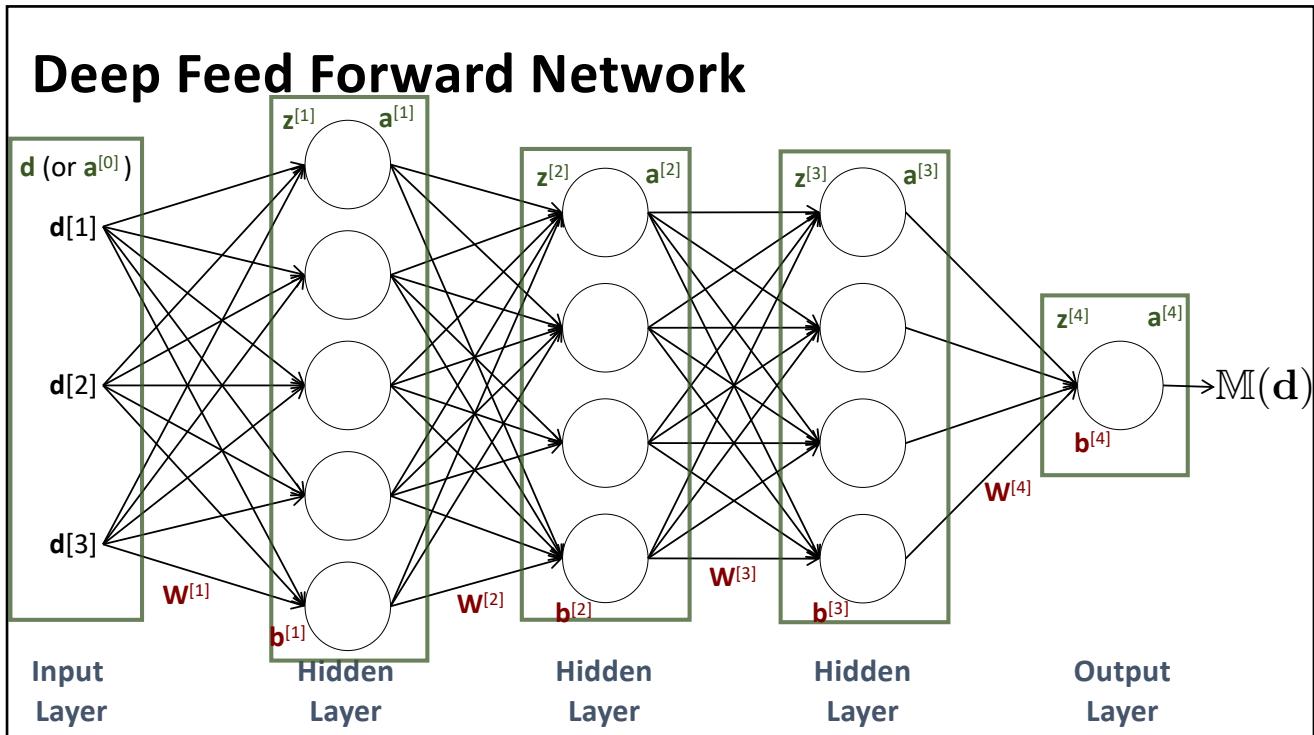


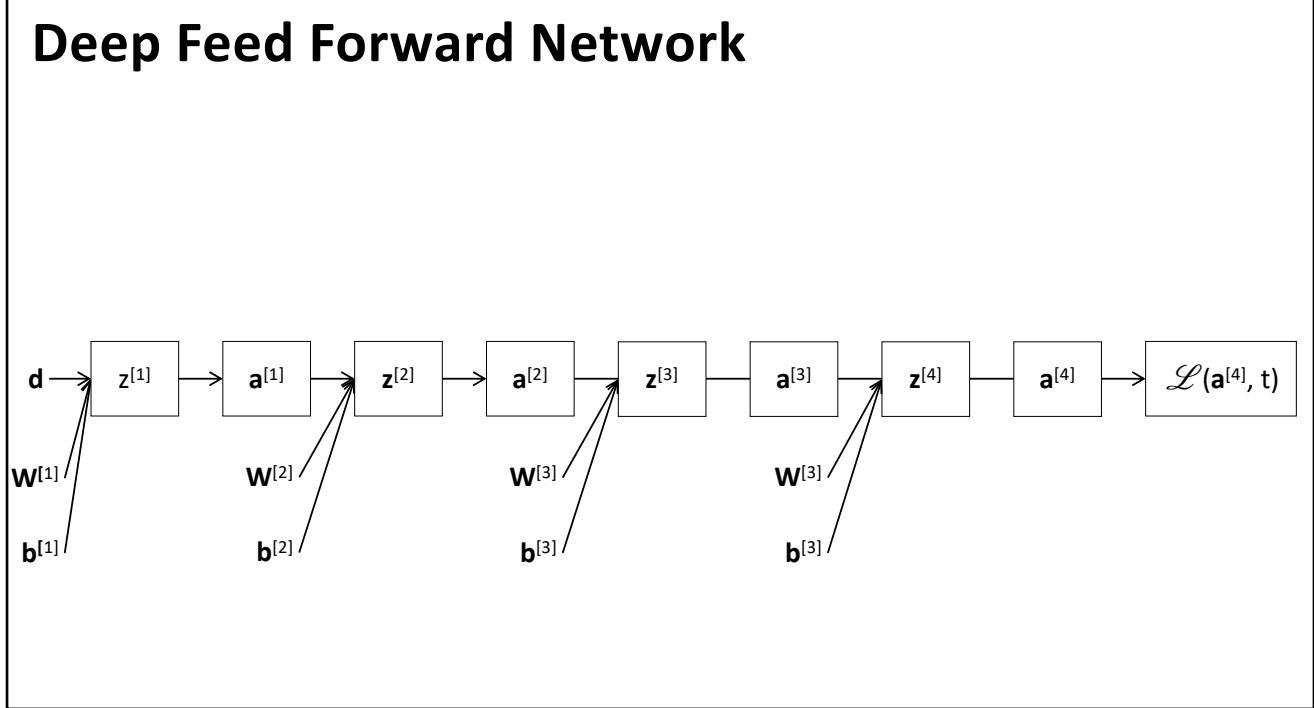
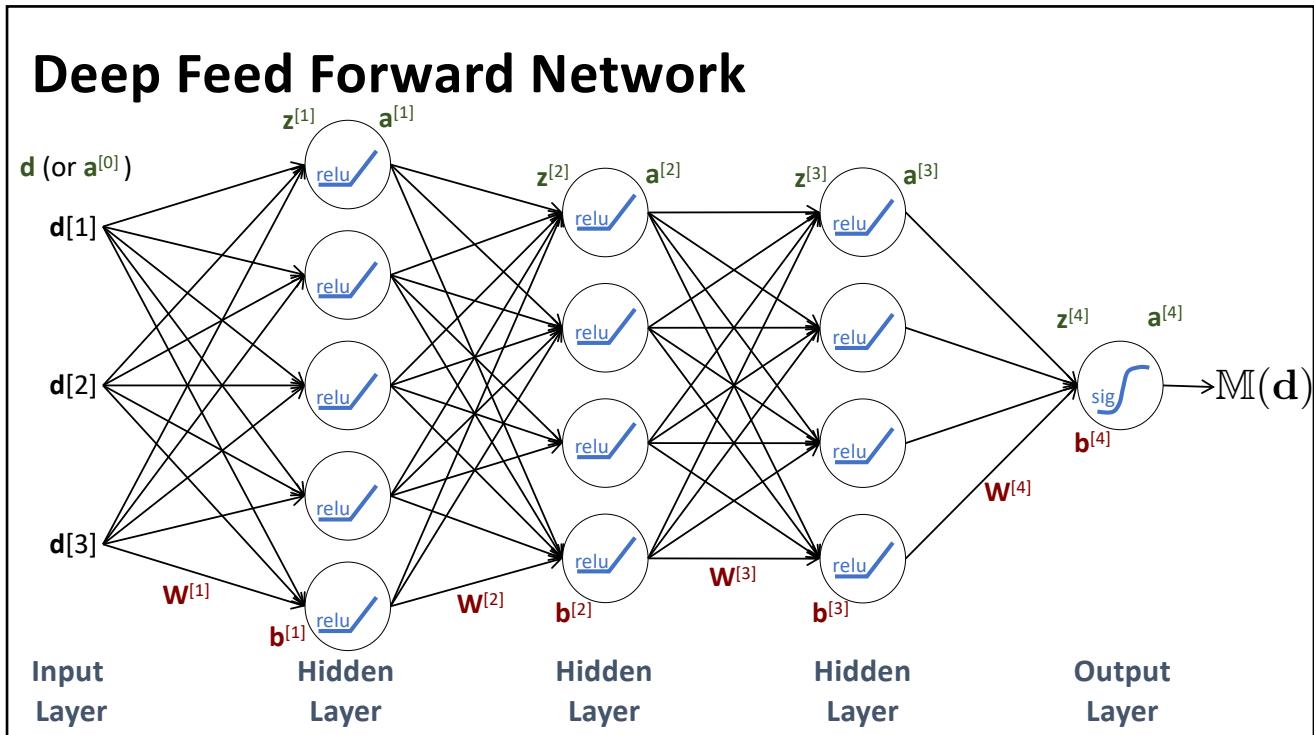
Deep Feed Forward Network



Deep Feed Forward Network







Deep Feed Forward Network Forward Propagation

Require: L, network depth

Require: $\mathbf{W}^{[i]}, i \in \{1, \dots, L\}$, weight matrices for each layer

Require: $\mathbf{b}^{[i]}, i \in \{1, \dots, L\}$, bias terms for each layer

Require: \mathbf{d} , input descriptive features

Require: t, the target feature

$$\mathbf{a}^{[0]} = \mathbf{d}$$

for l = 1 to L:

$$\mathbf{z}^{[l]} = \mathbf{W}^{[l]} \mathbf{a}^{[l-1]} + \mathbf{b}^{[l]}$$

$$\mathbf{a}^{[l]} = g^{[l]}(\mathbf{z}^{[l]})$$

$$\mathbb{M}(\mathbf{d}) = \mathbf{a}^{[L]}$$

Calculate $\mathcal{L}(\mathbb{M}(\mathbf{d}), t)$

Deep Feed Forward Network Backward Propagation

Require a forward pass through the network

Require: L, network depth

Require: $\mathbf{W}^{[i]}, i \in \{1, \dots, L\}$, weight matrices for each layer

Require: $\mathbf{b}^{[i]}, i \in \{1, \dots, L\}$, bias terms for each layer

Require: t, the target feature

Calculate $\mathbf{da}^{[L]}$ # The derivative of the loss function

for l = L to 1:

$$dz^{[l]} = da^{[l]} * g^{[l]}'(z^{[l]})$$

$$d\mathbf{W}^{[l]} = dz^{[l]} \mathbf{a}^{[l-1]T}$$

$$d\mathbf{b}^{[l]} = dz^{[l]}$$

$$da^{[l-1]} = \mathbf{W}^{[l]T} dz^{[l]}$$

Deep Feed Forward Network

Backward Propagation

Require a forward pass through the network

Require: L, network depth

Require: $\mathbf{W}^{[i]}$, $i \in \{1, \dots, L\}$, weight matrices for each layer

Require: $\mathbf{b}^{[i]}$, $i \in \{1, \dots, L\}$, bias terms for each layer

Require: t, the target feature

$da^{[L]} = -t/a^{[L]} + (1-t)/(1-a^{[L]})$ # The derivative of log loss

for $l = L$ to 1:

$$dz^{[l]} = da^{[l]} * g^{[l]}'(\mathbf{z}^{[l]})$$

$$d\mathbf{W}^{[l]} = dz^{[l]} \mathbf{a}^{[l-1]T}$$

$$d\mathbf{b}^{[l]} = dz^{[l]}$$

$$da^{[l-1]} = \mathbf{W}^{[l]T} dz^{[l]}$$

GRADIENT DESCENT

Gradient Descent Algorithm

Require: $\mathbb{M}_{\mathbf{W}, \mathbf{b}}$, a network upon which to perform gradient descent

Require: \mathcal{D} , a training dataset of m training instances $\{(\mathbf{d}_1, t_1), (\mathbf{d}_2, t_2), \dots, (\mathbf{d}_m, t_m)\}$

Require: $\mathcal{L}(\mathbb{M}(\mathbf{d}), t)$, a loss function to calculate the loss of a prediction

Require: $J(\mathbb{M}(\mathcal{D}_d), \mathcal{D}_t)$, a cost function to aggregate losses

Gradient Descent Algorithm (Stochastic)

Choose random weights

Until convergence

- Shuffle all training examples
- For each training instance
 - Perform a forward pass through the network
 - Perform a backward pass through the network
 - Update weights and biases using update rule

Gradient Descent Algorithm (Stochastic) Update Rule

We apply the same simple update rule at every layer

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

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

Gradient Descent Algorithm (Batch)

Choose random weights

Until convergence

- Set all gradient sums to 0
- For each training instance
 - Perform a forward pass through the network
 - Perform a backward pass through the network
 - Update gradient sum for each weight and bias term
- Update weights and biases using update rule

Cost Function

$$J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\mathbb{M}(\mathbf{d}_i), t_i)$$

$$\frac{\partial J(\mathbf{w}, b)}{\partial \mathbf{w}[j]} = \frac{1}{m} \sum_{i=1}^m \frac{\partial \mathcal{L}(\mathbb{M}(\mathbf{d}_i), t_i)}{\partial \mathbf{w}[j]}$$

Gradient Descent Algorithm (Batch) Update Rule

We apply the same simple update rule at every layer

$$\mathbf{W}^{[i]} = \mathbf{W}^{[i]} - \alpha \frac{1}{m} \sum_{j=1}^m \mathbf{d}\mathbf{W}^{[i]}_j$$

$$\mathbf{b}^{[i]} = \mathbf{b}^{[i]} - \alpha \frac{1}{m} \sum_{j=1}^m \mathbf{d}\mathbf{b}^{[i]}_j$$

Gradient Descent Algorithm (Mini Batch)

Choose random weights

Until convergence

- Divide the training dataset into mini-batches of size s
- For each mini-batch \mathcal{D}_{mb}
 - Set all gradient sums to 0
 - For each training instance in \mathcal{D}_{mb}
 - Perform a forward pass through the network
 - Perform a backward pass through the network
 - Update gradient sum for each weight and bias term
 - Update weights and biases using update rule

Gradient Descent Algorithm (Batch)

Update Rule

We apply the same simple update rule at every layer

$$\mathbf{W}^{[i]} = \mathbf{W}^{[i]} - \alpha \frac{1}{S} \sum_{j=1}^S \mathbf{dW}^{[i]}_j$$

$$\mathbf{b}^{[i]} = \mathbf{b}^{[i]} - \alpha \frac{1}{S} \sum_{j=1}^S \mathbf{db}^{[i]}_j$$

Different Variants of Gradient Descent

Stochastic gradient descent is easy to implement and can result in fast learning, **but** is computationally expensive (no opportunity to speed up) and can result in a noisy gradient signal

Batch gradient descent is computationally efficient and can result in a stable gradient signal, **but** requires gradient accumulation, can result in premature convergence can require loading large datasets into memory and can become slow

Machine Learning Mastery
<https://machinelearningmastery.com/easy-introduction-mini-batch-gradient-descent-configure-batch-size/>

Different Variants of Gradient Descent

Mini-batch gradient descent is relatively computationally efficient, does not require full datasets to be loaded into memory, and can result in a stable gradient signal, **but** requires gradient accumulation, and introduces a new hyper-parameter - mini-batch size

Machine Learning Mastery
<https://machinelearningmastery.com/easy-introduction-mini-batch-gradient-descent-configure-batch-size/>

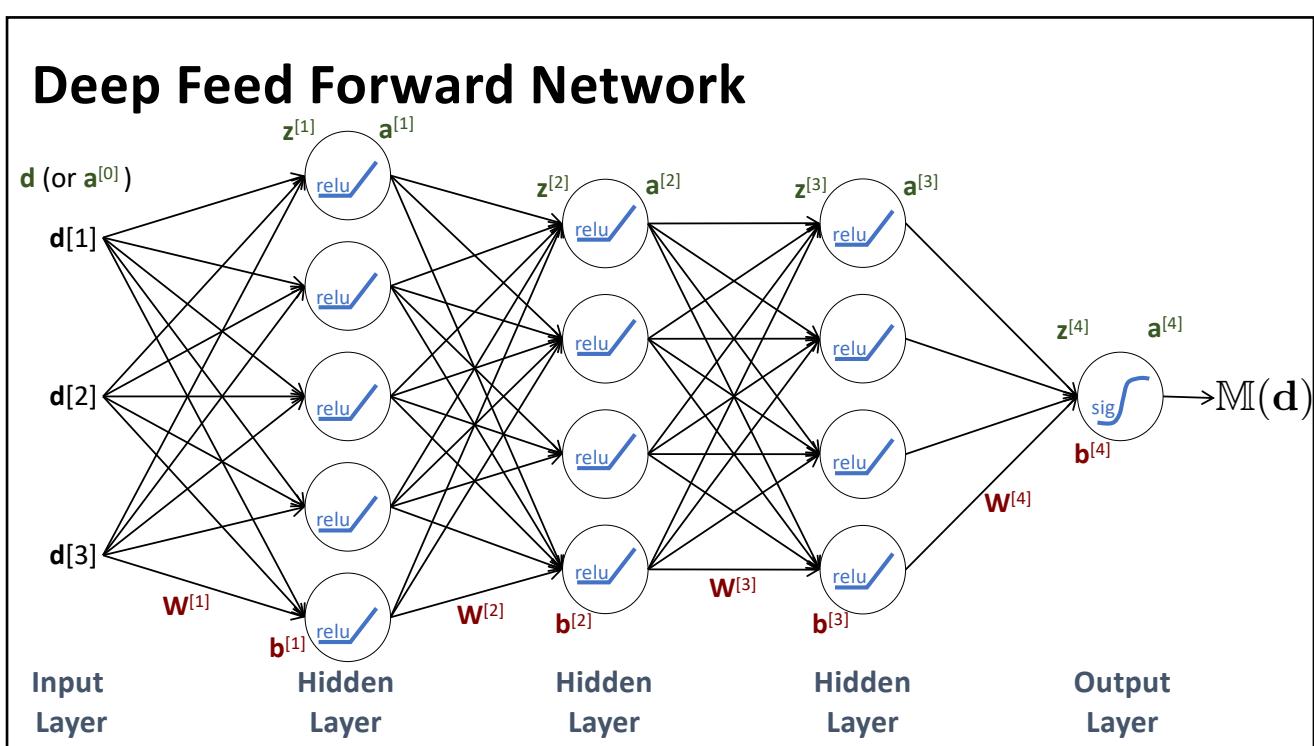
Different Variants of Gradient Descent

Most modern implementations use mini-batch gradient descent

Batch sizes are typically chosen to take advantage of computational infrastructure - e.g. 32, 64, 128,

...

OTHER LOSS FUNCTIONS



Multinomial Classification

Our combination of the sigmoid activation function and log loss loss function works fine for binary classification problems

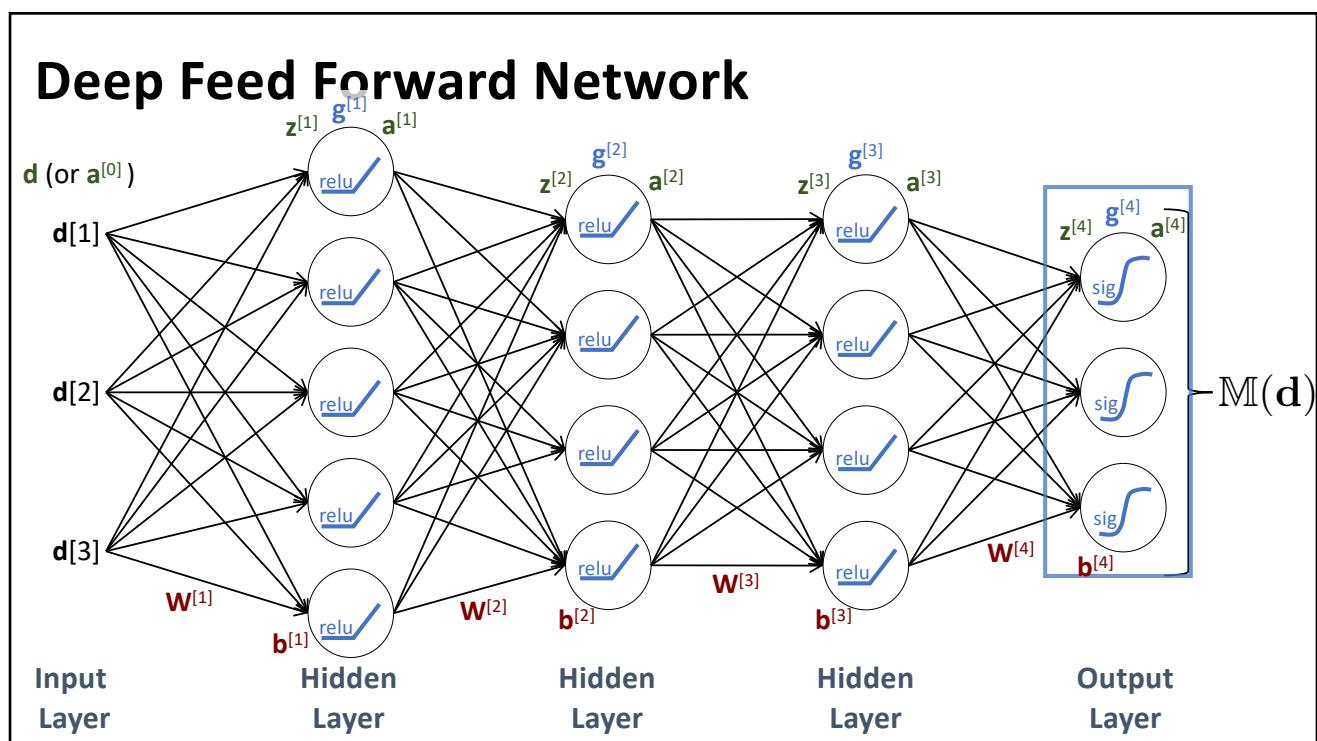
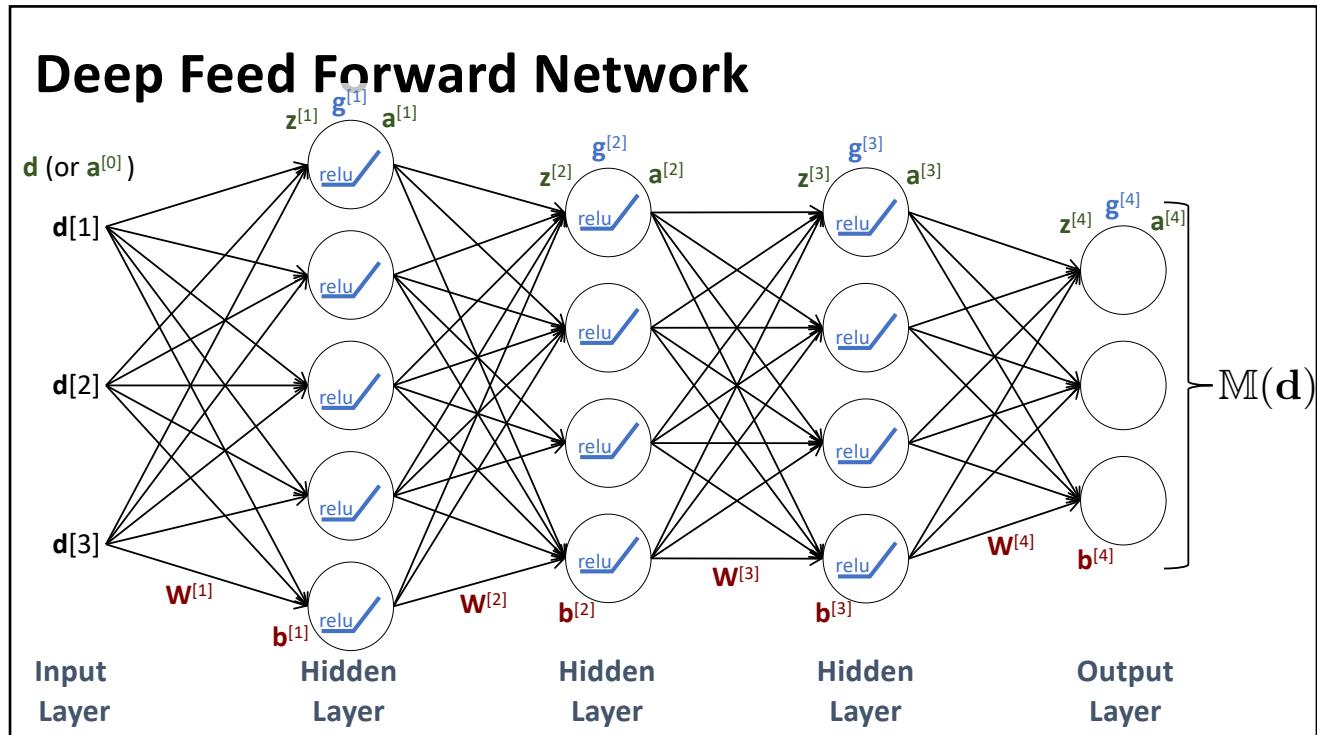
For multinomial classification problems we need to change things

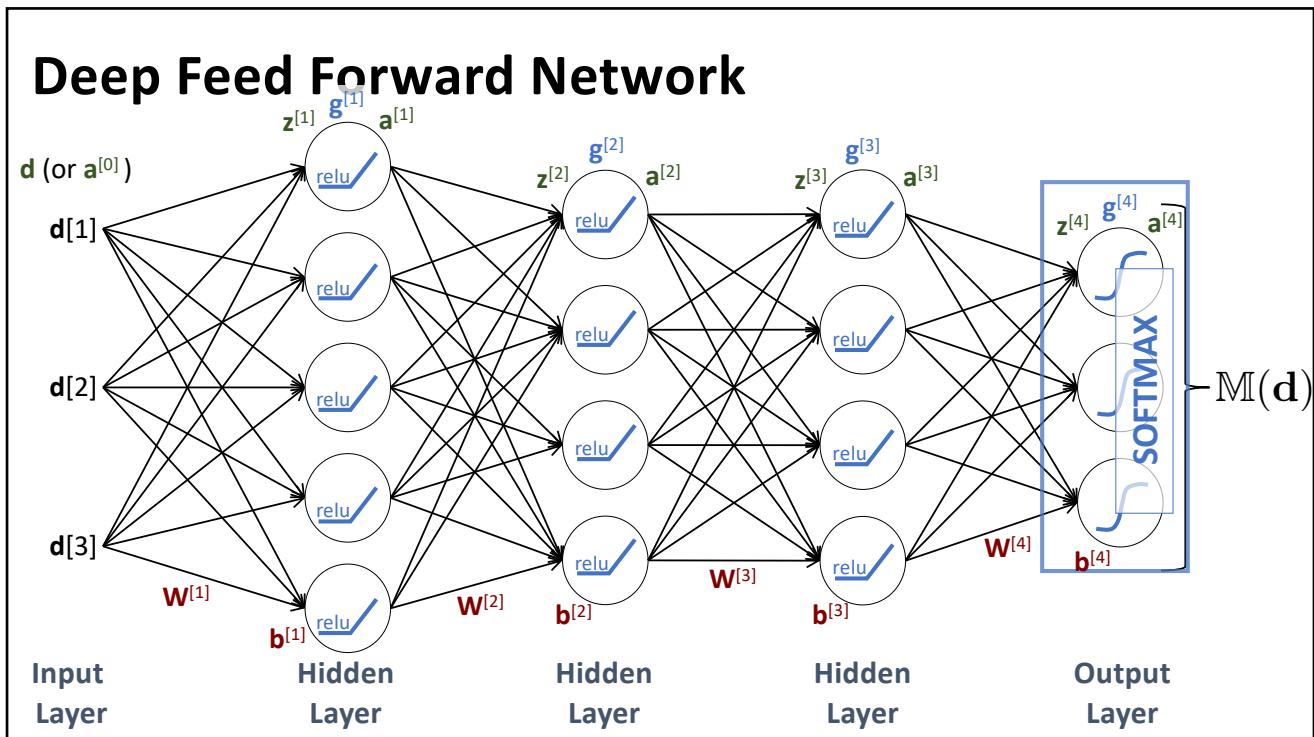
Multinomial Classification

	Descriptive Features					Target Feature
d_1	---	---	---	---	---	shirt
d_2	---	---	---	---	---	sneaker
			●			
			●			
			●			
d_{m-2}	---	---	---	---	---	trousers
d_{m-1}	---	---	---	---	---	shirt
d_m	---	---	---	---	---	sneaker

Multinomial Classification

	Descriptive Features	Target Feature		
		(shirt)	(sneaker)	(trousers)
d_1	---	1	0	0
d_2	---	0	1	0
d_{m-2}	---	0	0	1
d_{m-1}	---	1	0	0
d_m	---	0	1	0





Activation Function: Softmax

The **softmax** (or normalized exponential) function is a generalization of the sigmoid function that *squashes* a vector of arbitrary real values a vector real values in the range $[0, 1]$ that add up to 1

$$\text{softmax}(\mathbf{z}[i]) = \frac{e^{\mathbf{z}[i]}}{\sum_k e^{\mathbf{z}[k]}}$$

Activation Function: Softmax Example

Suppose that the output of our network is

$$\mathbf{z} = [5, -1, 2]$$

then

$$e^{\mathbf{z}} =$$

Activation Function: Softmax Example

Suppose that the output of our network is

$$\mathbf{z} = [5, -1, 2]$$

then

$$e^{\mathbf{z}} = [148.41, 0.34, 7.34]$$

then

$$\text{sum}(e^{\mathbf{z}}) =$$

Activation Function: Softmax Example

Suppose that the output of our network is

$$\mathbf{z} = [5, -1, 2]$$

then

$$e^{\mathbf{z}} = [148.41, 0.34, 7.34]$$

then

$$\text{sum}(e^{\mathbf{z}}) = 156.17$$

and

$$\text{softmax}[\mathbf{z}] =$$

Activation Function: Softmax Example

Suppose that the output of our network is

$$\mathbf{z} = [5, -1, 2]$$

then

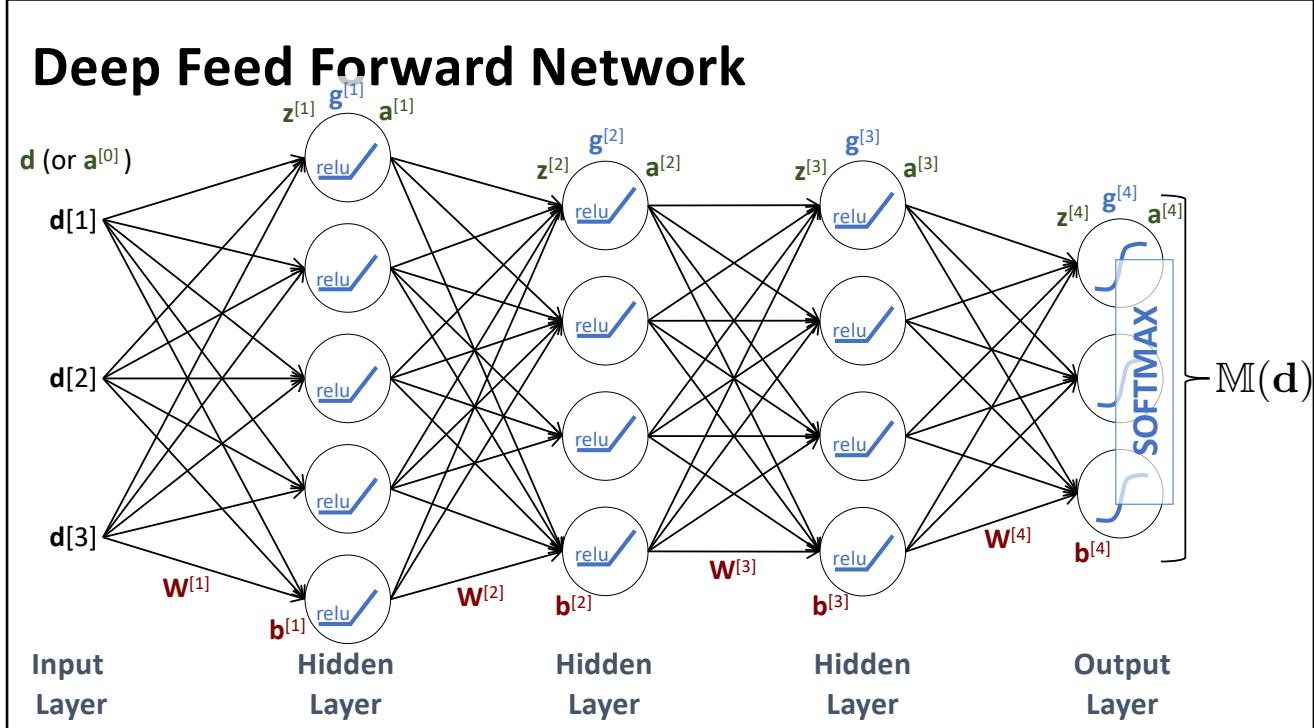
$$e^{\mathbf{z}} = [148.41, 0.34, 7.34]$$

then

$$\text{sum}(e^{\mathbf{z}}) = 156.17$$

and

$$\text{softmax}[\mathbf{z}] = [0.950, 0.002, 0.047]$$



Loss Function: Cross Entropy

Using multiple outputs requires a new loss function

The most common loss function used for this scenario is **cross entropy**

$$crossent(\mathbb{M}(\mathbf{d}), \mathbf{t}) = - \sum_{k=1}^l \mathbf{t}[k] \log(\mathbb{M}(\mathbf{d})[k])$$

Loss Function: Cross Entropy

Log loss is just a special case of cross entropy for two target feature levels

$$\begin{aligned}\mathcal{L}(\mathbb{M}(\mathbf{d}), t) = \\ -((t \times \log(\mathbb{M}(\mathbf{d}))) + (1 - t) \times \log(1 - \mathbb{M}(\mathbf{d}))))\end{aligned}$$

$$\begin{aligned}\mathcal{L}(\mathbb{M}(\mathbf{d}), \mathbf{t}) = \\ -((\mathbf{t}[1] \times \log(\mathbb{M}(\mathbf{d}[1]))) + (\mathbf{t}[2] \times \log(\mathbb{M}(\mathbf{d})[2])))\end{aligned}$$

$$crossent(\mathbb{M}(\mathbf{d}), \mathbf{t}) = - \sum_{k=1}^2 \mathbf{t}[k] \log(\mathbb{M}(\mathbf{d})[k])$$

Activation Function: Softmax Example

Suppose that the output of our network after softmax is

$$\mathbb{M}(\mathbf{d}) = [0.950, 0.002, 0.047]$$

and the target feature vector, \mathbf{t} , is

$$\mathbf{t} = [1, 0, 0]$$

then

$$\begin{aligned}crossent(\mathbb{M}(\mathbf{d}), \mathbf{t}) \\ =\end{aligned}$$

Activation Function: Softmax Example

Suppose that the output of our network after softmax is

$$\mathbb{M}(\mathbf{d}) = [0.950, 0.002, 0.047]$$

and the target feature vector, \mathbf{t} , is

$$\mathbf{t} = [1, 0, 0]$$

then

$$\begin{aligned} & \text{crossent}(\mathbb{M}(\mathbf{d}), \mathbf{t}) \\ &= \log(0.950) + 0 \log(0.002) + 0 \\ & \quad \log(0.047) \\ &= -(-0.051) = 0.051 \end{aligned}$$

Loss Function: Cross Entropy Derivative

Conveniently the derivative of cross entropy loss with respect to the output of a unit in the final layer of a network is:

$$\frac{\partial \text{crossent}(\mathbb{M}(\mathbf{d}), \mathbf{t})}{\partial \mathbf{a}[i]} = \mathbb{M}(\mathbf{d})[i] - \mathbf{t}[i]$$

Cost Function: Cross Entropy

And just like we did before for batch or min-batch gradient descent we can sum the cross entropy losses across a dataset to generate a cost function

$$J_{crossent}(\mathcal{D}, \mathbb{M}) = \frac{1}{m} \sum_{i=1}^m crossent(\mathbb{M}(\mathbf{d}_i), \mathbf{t}_i)$$

Cost Function: Cross Entropy

And also like before for batch or min-batch gradient descent we can sum the cross entropy loss derivatives across a dataset to generate the derivative of the cost function

$$\frac{\partial J_{crossent}(\mathcal{D}, \mathbb{M})}{\partial \mathbf{w}[j]} = \frac{1}{m} \sum_{i=1}^m \frac{\partial \mathcal{L}(\mathbb{M}(\mathbf{d}_i), \mathbf{t}_i)}{\partial \mathbf{w}[j]}$$

Deep Feed Forward Network Backward Propagation

Require a forward pass through the network

Require: L, network depth

Require: $\mathbf{W}^{[i]}$, $i \in \{1, \dots, L\}$, weight matrices for each layer

Require: $\mathbf{b}^{[i]}$, $i \in \{1, \dots, L\}$, bias terms for each layer

Require: t, the target feature

Calculate $\mathbf{da}^{[L]}$ # The derivative of the loss function
for $l = L$ to 1:

$$dz^{[l]} = da^{[l]} * g^{[l]}'(\mathbf{z}^{[l]})$$

$$d\mathbf{W}^{[l]} = dz^{[l]} \mathbf{a}^{[l-1]T}$$

$$d\mathbf{b}^{[l]} = dz^{[l]}$$

$$da^{[l-1]} = \mathbf{W}^{[l]T} dz^{[l]}$$

Deep Feed Forward Network Backward Propagation

Require a forward pass through the network

Require: L, network depth

Require: $\mathbf{W}^{[i]}$, $i \in \{1, \dots, L\}$, weight matrices for each layer

Require: $\mathbf{b}^{[i]}$, $i \in \{1, \dots, L\}$, bias terms for each layer

Require: t, the target feature

$da^{[L]} = t/a^{[L]}$ # The derivative of the loss function
for $l = L$ to 1:

$$dz^{[l]} = da^{[l]} * g^{[l]}'(\mathbf{z}^{[l]})$$

$$d\mathbf{W}^{[l]} = dz^{[l]} \mathbf{a}^{[l-1]T}$$

$$d\mathbf{b}^{[l]} = dz^{[l]}$$

$$da^{[l-1]} = \mathbf{W}^{[l]T} dz^{[l]}$$

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

