



# Intro to Neural Networks

Foundations of Deep Learning

BOSTON UNIVERSITY  
MACHINE INTELLIGENCE  
COMMUNITY

Presenter: Julius  
09/11/2019

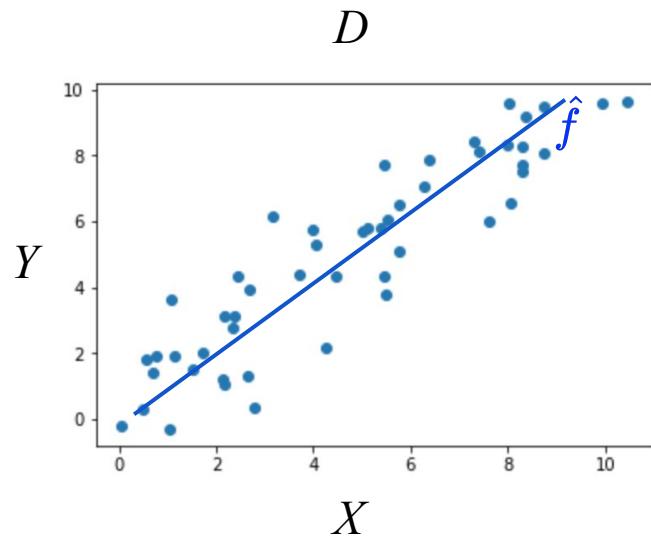


# Building up to Neural Networks

Regression

# What is learning?

The approximation of some unknown function  $f$  based on some data  $D$ .



$$f : X \rightarrow Y$$

$$\hat{f} = \theta_0 + \theta_1 x$$

How do we set the parameters?  
How do we know what assumptions to make?

# Dataset

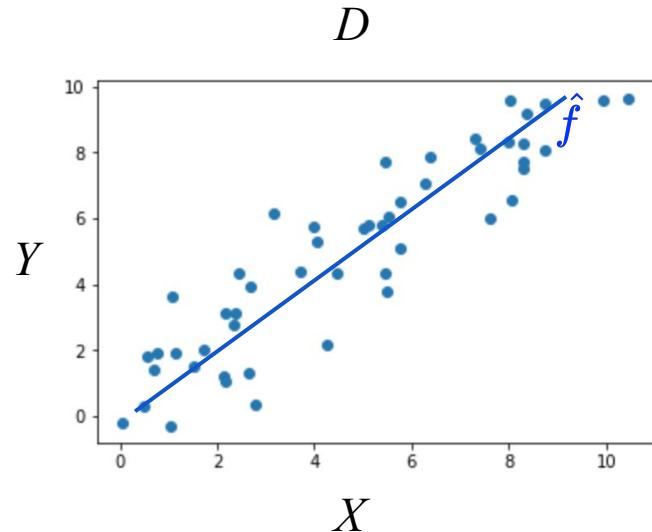
$$X_i = \begin{bmatrix} feature_1 \\ feature_2 \\ \vdots \\ feature_n \end{bmatrix}$$

$$X = \begin{array}{|c|c|c|}\hline & m_{examples} & \\\hline & \hline & \\\hline\end{array} \quad \left. \right\} n_{features}$$

$$Y_i = label_1$$

$$Y = \begin{array}{|c|c|c|}\hline & m_{examples} & \\\hline & \hline & \\\hline\end{array}$$

# Linear Regression



One dimensional input

$$y = mx + b$$

Multi dimensional input

$$y = w_1 x_1 + \dots + w_n x_n + b$$

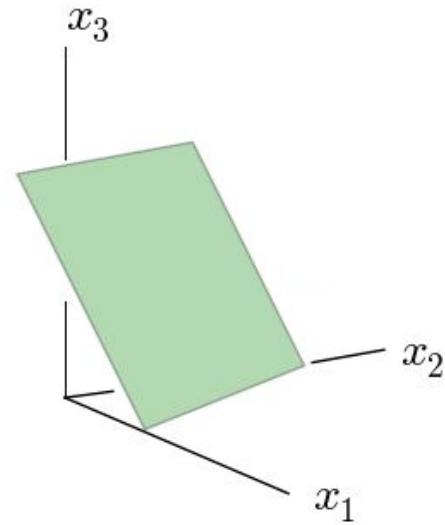
# Linear Regression example

**Say you want to predict house prices given some features about the house**

$$x_1 = 5(\text{bedrooms}), x_2 = 3(\text{bathrooms}), x_3 = 3000(\text{squarefeet})$$

**Then the house price could be:**

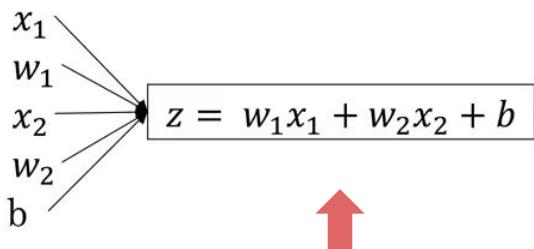
$$z = 5w_1 + 3w_2 + 3000w_3 + b$$



# Logistic Regression (1)

1. Has a linear combination / made up of weights and a bias

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$



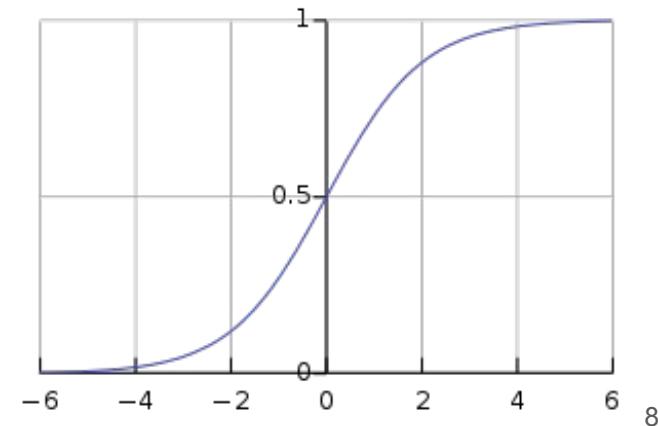
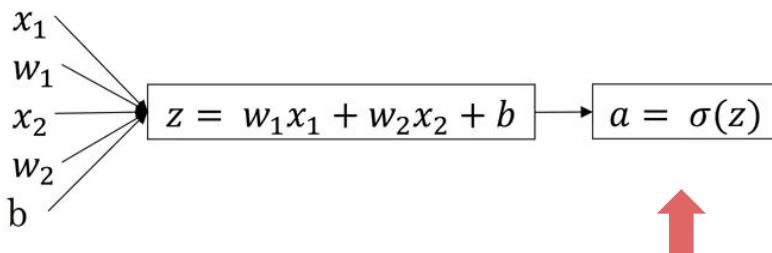
# Logistic Regression (2)

1. Has a linear combination / made up of weights and a bias

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

2. Is fed through an activation function (such as a sigmoid function)

$$a = \hat{y} = \sigma(z)$$



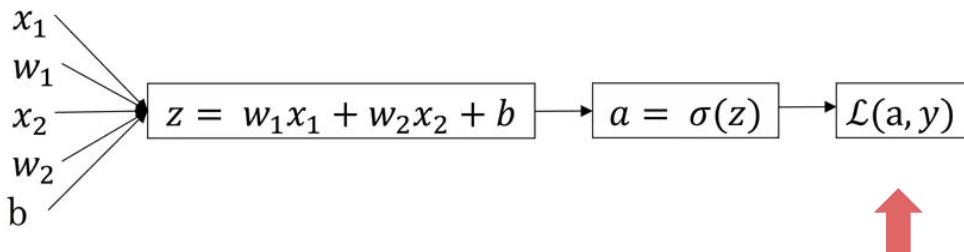
# Logistic Regression (3)

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

$$a = \hat{y} = \sigma(z)$$

3. The output is used to measure the difference from the actual data.

$$L(a, y) = y - a$$



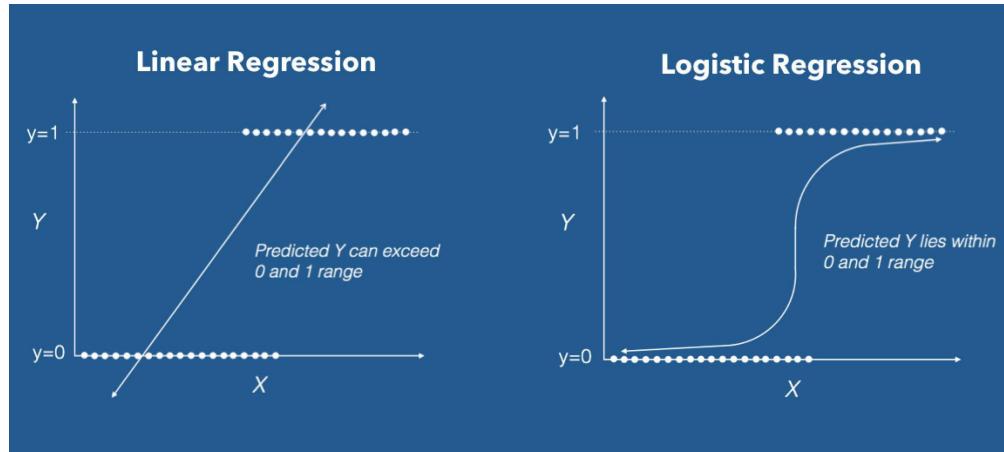
# Logistic Regression (example)

Now say we want to determine if we can afford the house or not.

In this case we would use a sigmoid function to output values between 0 or 1.

This can be a metric of how affordable the house is.

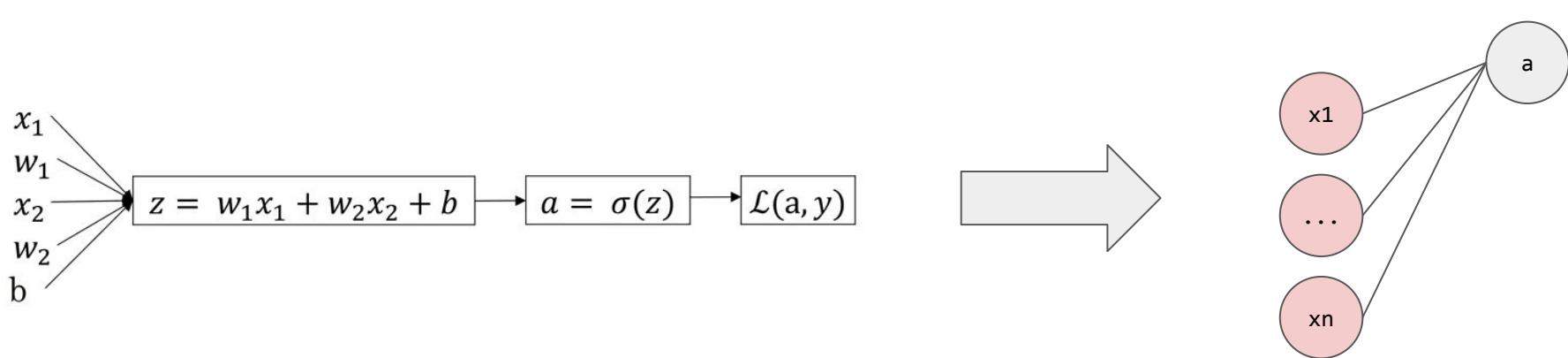
$$a = y = \sigma(z)$$



# Logistic Regression to Neural Networks (1)

Single logistic regression

$$a = \sigma(w_1x_2 + \dots + w_nx_n + b)$$



Expanded representation

Compact representation

# Logistic Regression to Neural Networks (2)

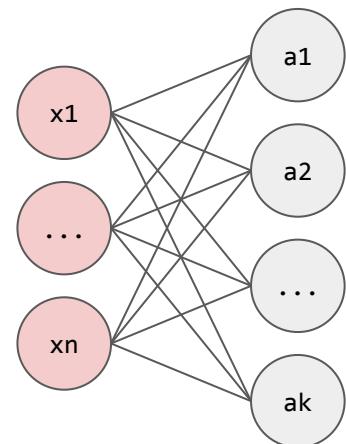
Multiple logistic regressions

$$a_1 = \sigma(w_1^{(1)}x_1 + \cdots + w_n^{(1)}x_n + b^{(1)})$$

$$a_2 = \sigma(w_1^{(2)}x_1 + \cdots + w_n^{(2)}x_n + b^{(2)})$$

⋮

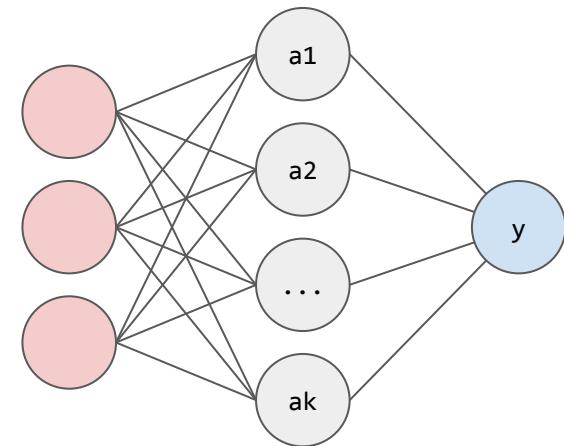
$$a_k = \sigma(w_1^{(k)}x_1 + \cdots + w_n^{(k)}x_n + b^{(k)})$$



# Logistic Regression to Neural Networks (2)

**Logistic regression on the outputs of other logistic regressions**

$$\hat{y} = \sigma(w_1^{(y)} a_1 + \dots + w_k^{(y)} a_k + b^{(y)})$$



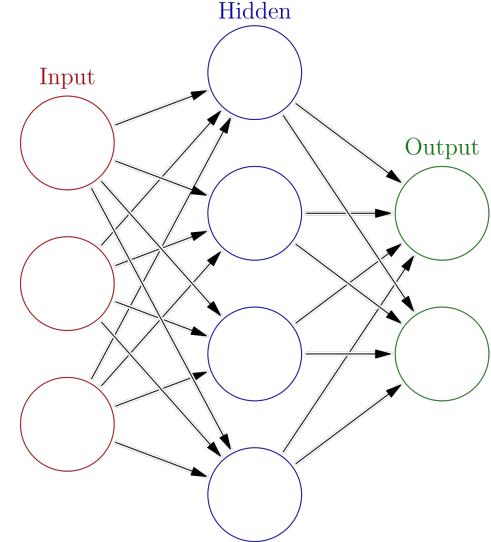


# Shallow Neural Networks

The foundation of new AI

# Neural Network

1. A neural network is a bunch of activations in layers
  - a. The middle layers are called a **hidden layer**
2. Each activation takes in inputs from the previous layer
  - a. The **hidden layer** is built from the **input layer**
  - b. The **output layer** is built from the **hidden layer**
3. Adding more layers means we can learn more complicated functions



# Activations (Neuron)

## 1. Affine function

a.  $z = w_1 a_1^{[l-1]} + \cdots + w_n a_n^{[l-1]} + b$

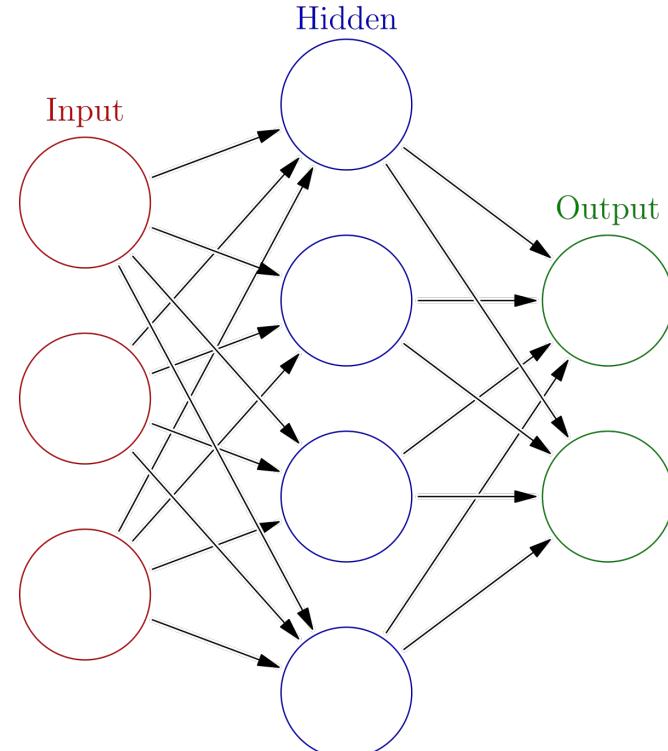
b. Gathers and weights features

## 2. Activation

a.  $a = g(z)$

i.  $g(x)$  is an activation function

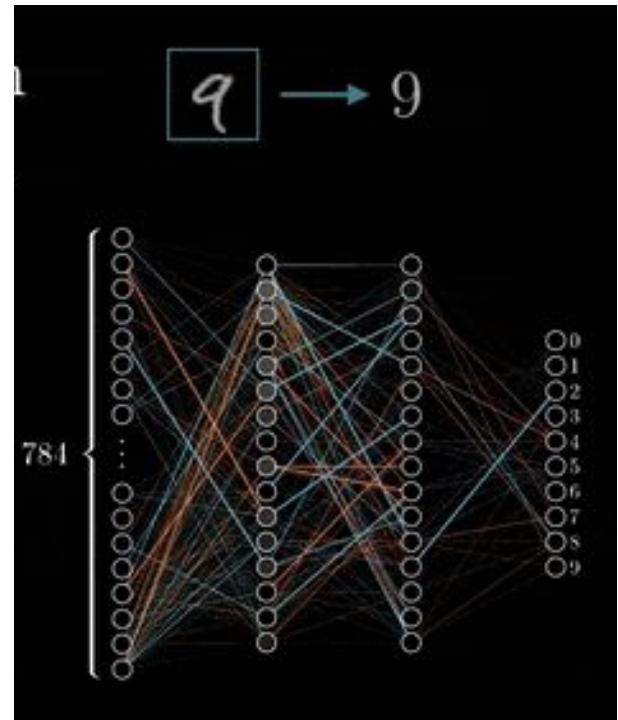
b. Acts as a significance threshold



Where are the weights in this diagram?

# Three steps to training a neural network

1. Forward propagation
  - a. Push example through the network to get a predicted output
2. Compute the cost
  - a. Calculate difference between predicted output and actual data
3. Backward propagation
  - a. Push back the derivative of the error and apply to each weight, such that next time it will result in a lower error



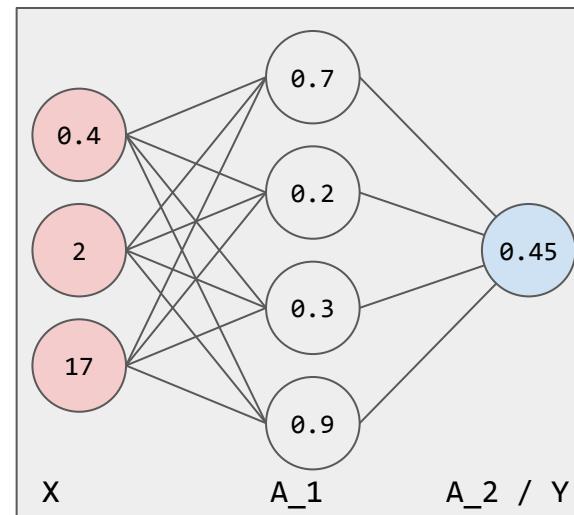
# Shallow Neural Networks (example)

**Say you want to predict if it is worth going to class or not**

$$X = \begin{bmatrix} \text{difficulty of material} \\ \text{number of assignments due} \\ \text{days until midterm} \end{bmatrix} \quad Y = [\text{probability of going to class}]$$

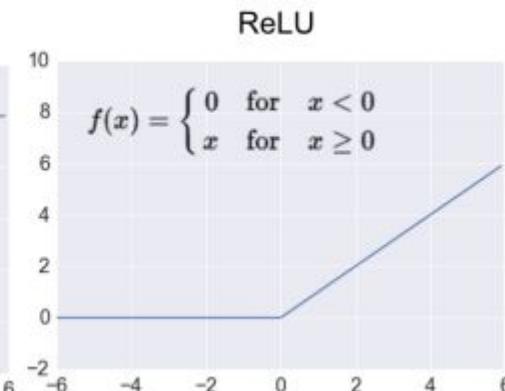
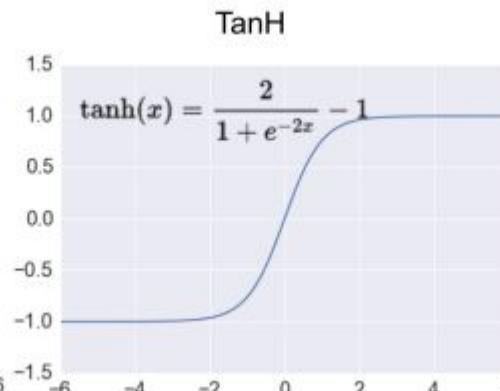
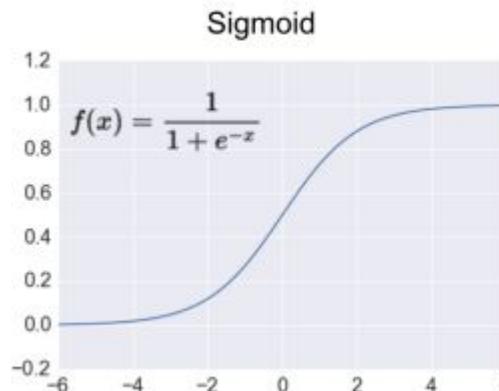
$$a_i^{[1]} = g(z) = g(w_1x_1 + w_2x_2 + w_3x_3 + b)$$

$$a^{[2]} = g(w_1a_1^{[1]} + w_2a_2^{[1]} + w_3a_3^{[1]} + w_4a_4^{[1]} + b)$$



# Activation Functions

1. Sigmoid: output is between 0,1
2. Tanh: output is between -1,1
3. ReLU: output is positive real numbers



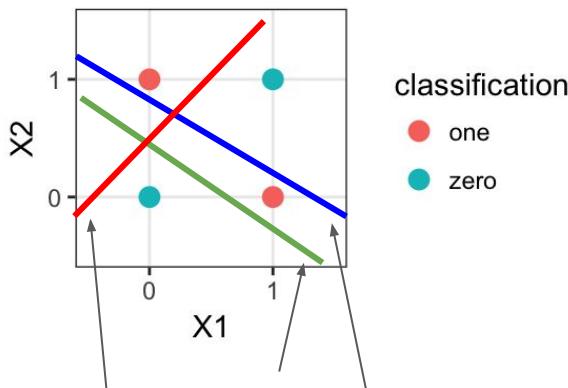
# Why do you need nonlinear activation functions?

## XOR Problem

*There is no way to correctly classify all inputs with a linear decision boundary*

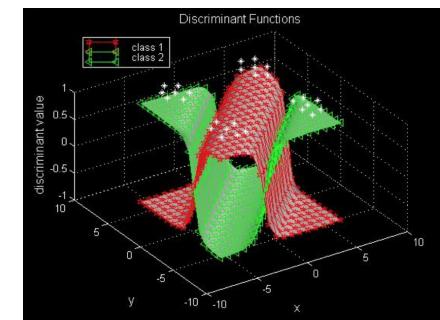
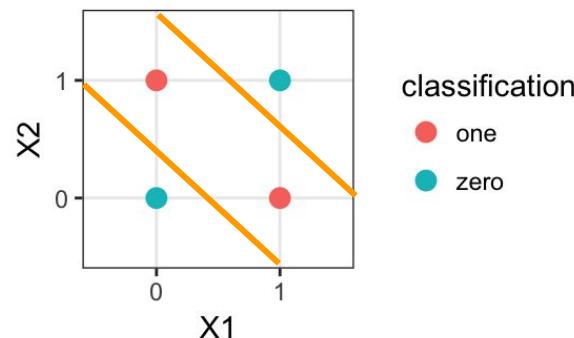
Linear Classifier

Classification of XOR



Two Layer Neural Network w/  
Sigmoid Activation

Classification of XOR



# Linear Activation Functions

If we removed the activation function from our model that can be called a linear activation function.

$$a_1 = w_1 x + b_1 \quad a_2 = w_2 a_1 + b_2$$

$$\begin{aligned} a_2 &= w_2(w_1 x + b_1) + b_2 \\ &= w_2 w_1 x + \underbrace{w_2 b_1 + b_2}_{\text{bias}} \\ &= wx + b \end{aligned}$$



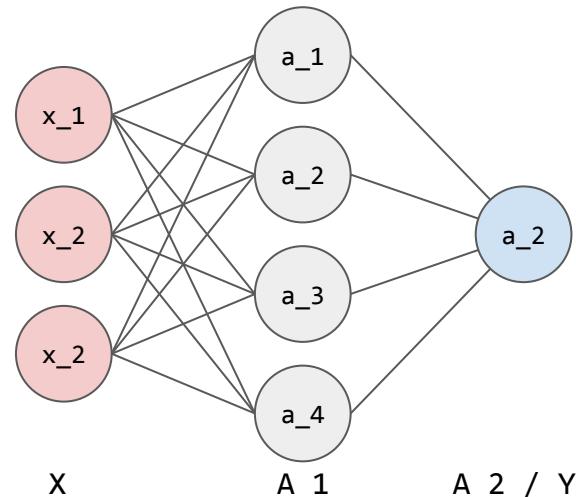
# Forward Propagation

1. We've seen neural nets with one pass through, however usually, we have multiple samples.
2. One way is to do this iteratively

```
for each sample i:  
    for each layer l:  
        for each activation j:
```

$$a_j^{[l]} = g^{[l]} \left( \sum_k w_{jk}^{[l]} a_k^{[l-1]} + b_j^{[l]} \right) = g^{[l]}(z_j^{[l]})$$

finally, use the final activations to compute the cost



# Computing the cost

In order to train our neural network, we need some way to tell us how far off its estimate was from the actual value.

We define the cost function,  $J(\hat{y}, y)$  as the sum of losses,  $\sum_{i=0}^m L(\hat{y}, y)$

- a. Loss = Error for a single training example
- b. Cost = Sum of all Losses



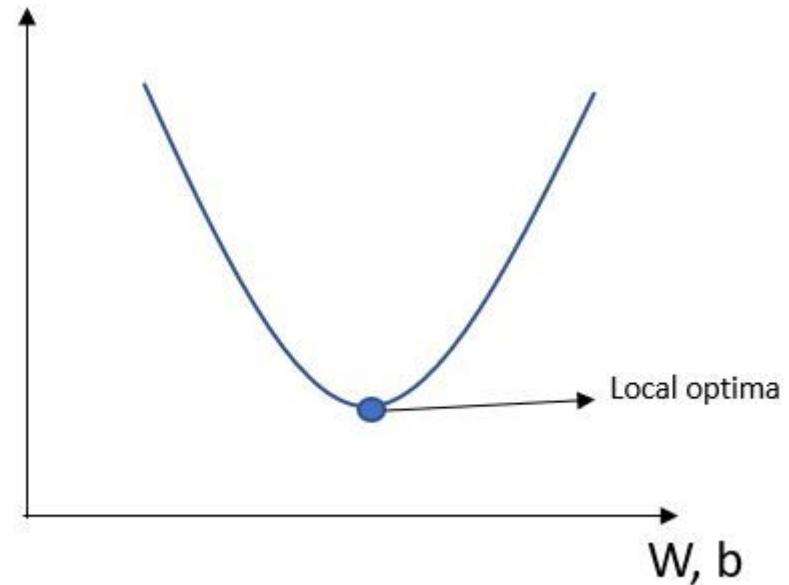
# Example cost functions

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

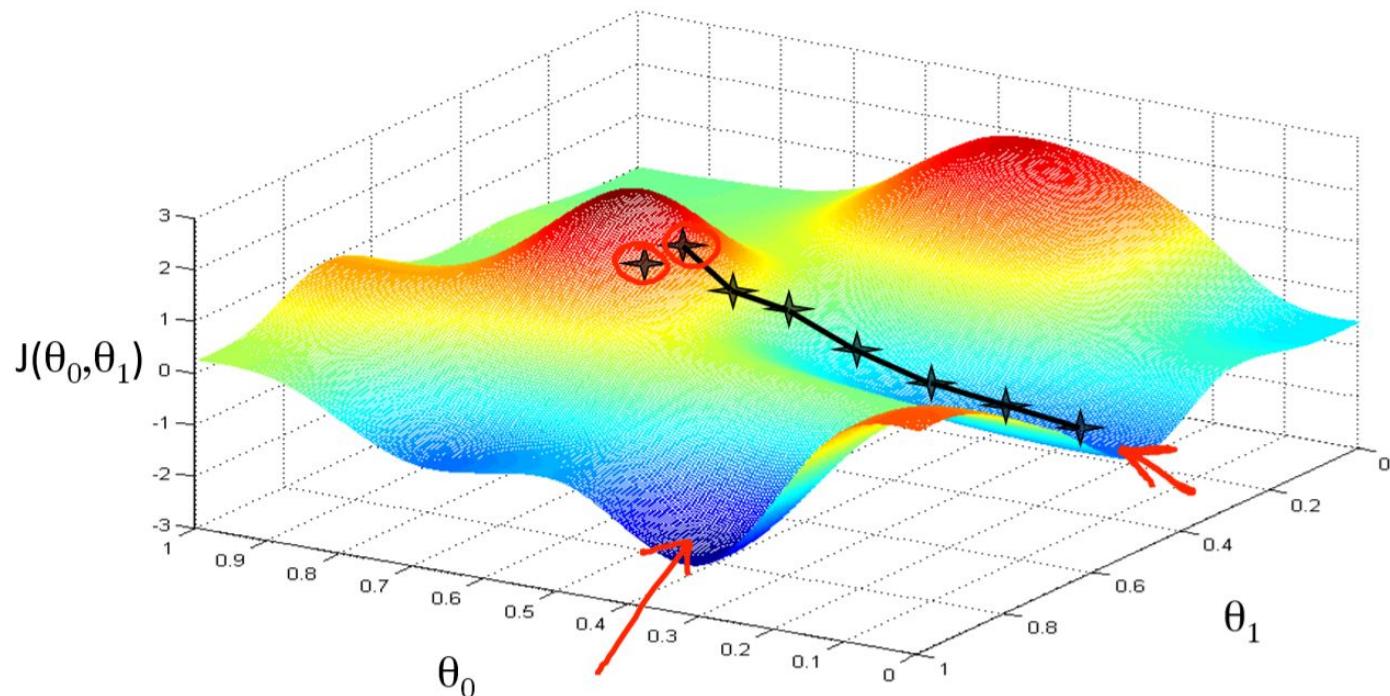
$$L_{\text{cross-entropy}}(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_i y_i \log(\hat{y}_i)$$

$$L(\hat{y}, y) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

Cost Function ( $J$ )

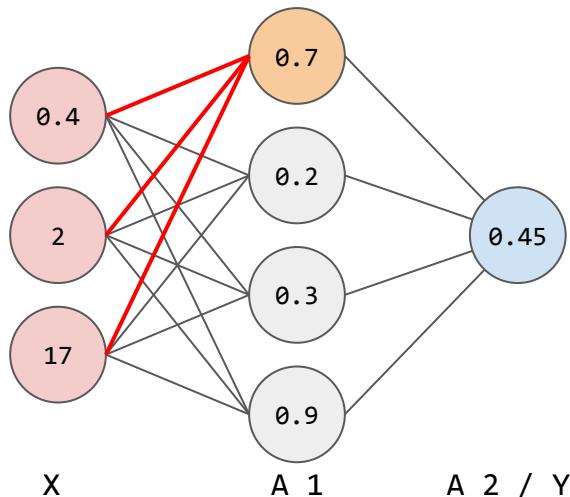


# Cost function for gradient descent

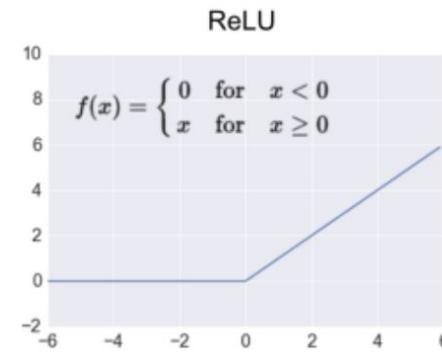


# Forward Propagation (example)

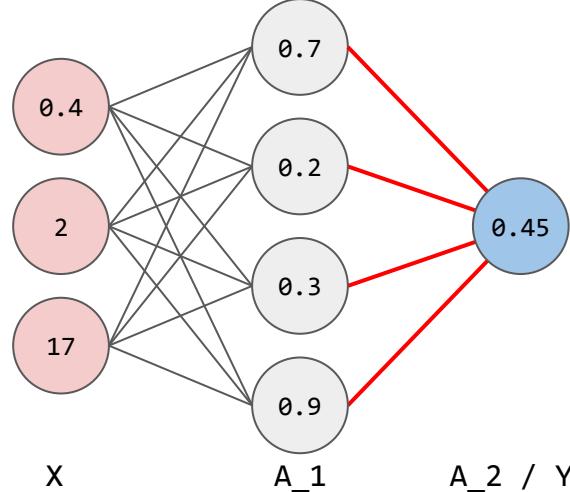
$$\begin{aligned}z_1^{[1]} &= w_1^{[1]}x_1 + w_2^{[1]}x_2 + w_3^{[1]}x_3 + b \\&= 0.5 * 0.4 + 0.1 \times 2 + 0.0058 \times 17 + 0.2 = 0.7\end{aligned}$$



$$a_1^{[1]} = g(0.7) = \text{ReLU}(0.7) = 0.7$$



# Forward Propagation (example)

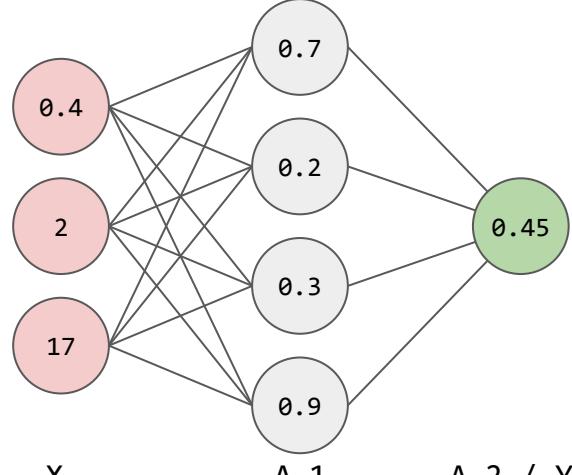


$$\begin{aligned}z_1^{[1]} &= w_1^{[1]}x_1 + w_2^{[1]}x_2 + w_3^{[1]}x_3 + b \\&= 0.5 * 0.4 + 0.1 \times 2 + 0.0058 \times 17 + 0.2 = 0.7\end{aligned}$$

$$a_1^{[1]} = g(0.7) = \text{ReLU}(0.7) = 0.7$$

$$z_1^{[2]} = w_1^{[2]}a_1^{[1]} + w_2^{[2]}a_2^{[1]} + w_3^{[2]}a_3^{[1]} + w_4^{[2]}a_4^{[1]} + b$$

# Forward Propagation (example)



$$\begin{aligned}z_1^{[1]} &= w_1^{[1]}x_1 + w_2^{[1]}x_2 + w_3^{[1]}x_3 + b \\&= 0.5 * 0.4 + 0.1 \times 2 + 0.0058 \times 17 + 0.2 = 0.7\end{aligned}$$

$$a_1^{[1]} = g(0.7) = \text{ReLU}(0.7) = 0.7$$

$$z_1^{[2]} = w_1^{[2]}a_1^{[1]} + w_2^{[2]}a_2^{[1]} + w_3^{[2]}a_3^{[1]} + w_4^{[2]}a_4^{[1]} + b$$

$$\hat{y} = a_1^{[2]} = g(z_1^{[2]}) = 0.45$$

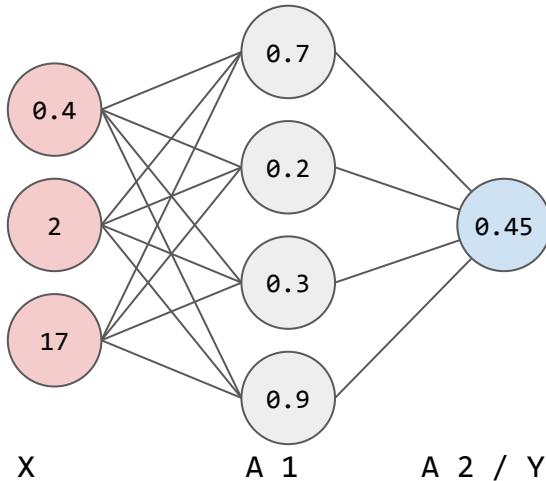
$$y = 1$$

$$L(\hat{y}, y) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

$$-(1) \log (0.45) - (1 - (1)) \log(1 - (0.45)) = 0.798508$$

# Back to the example

Say you want to predict if it is worth going to class or not



$$\hat{y} = a_1^{[2]} = g(z_1^{[2]}) = 0.45$$

```
if y > 0.5:  
    goto_class()  
else:  
    skip()
```

But it turns out midterm material was covered...

$$y = 1$$

Here is his mistake in numerical form

$$L(\hat{y}, y) = 0.798508$$



# Training a network with Calculus

Gradient descent

# Terminology for later - Parameters vs Hyperparameters

Parameters are values that are learned through training or backpropagation. (E.g.  $W$ ,  $b$ )

Hyperparameters are values that are set manually without a learning method. (E.g. the learning rate,  $\alpha$ )

# Cost Function: revisited

The cost function is the sum of losses for all training examples

- combined error for all training examples.

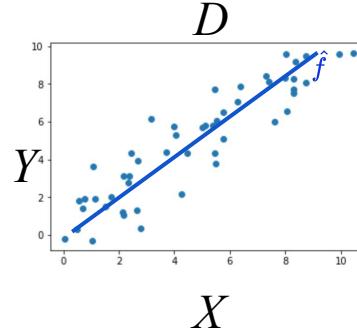
$$J(\hat{y}, y) = \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

This is the general form of writing the cost function

$$J(\theta) = \text{Cost Function}$$

# Learning the weights

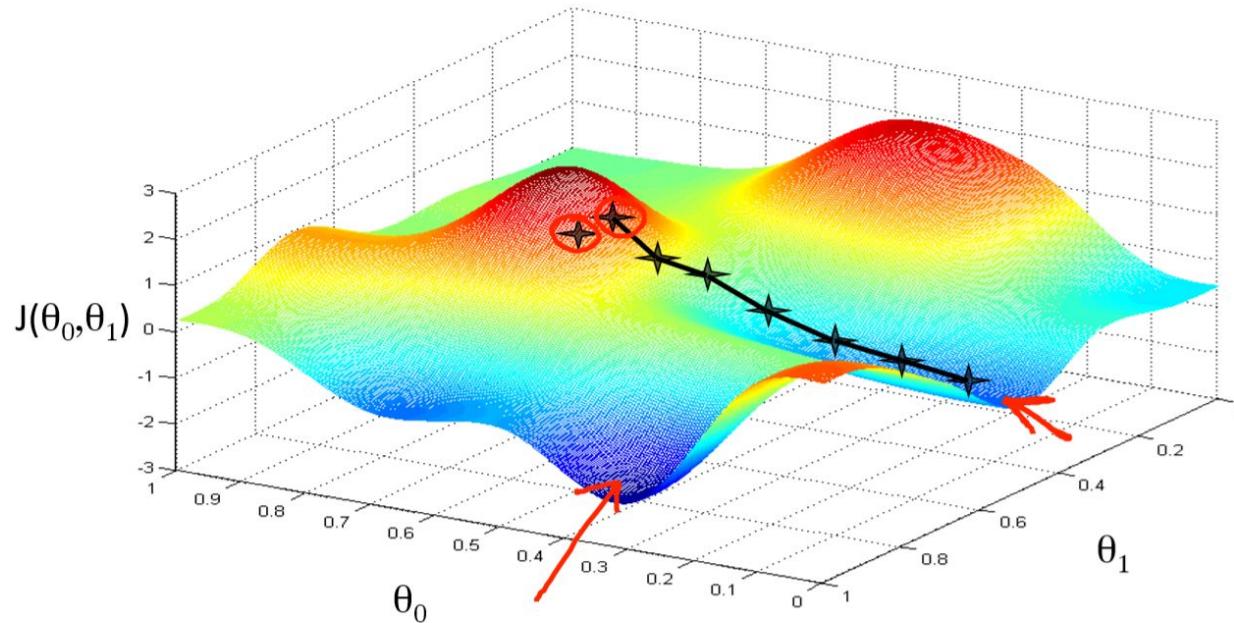
1. Goal: we want to change all of the weights so that our predictions fit the data better



2. Method: the error between our predicted values and actual values tells us how much we need to change each individual weight (increase or decrease)
  - a. By calculating the derivative from the cost function we get the numerical value of change.
  - b. The derivative is equivalent to the slope at that point, and tells us where the error is minimized

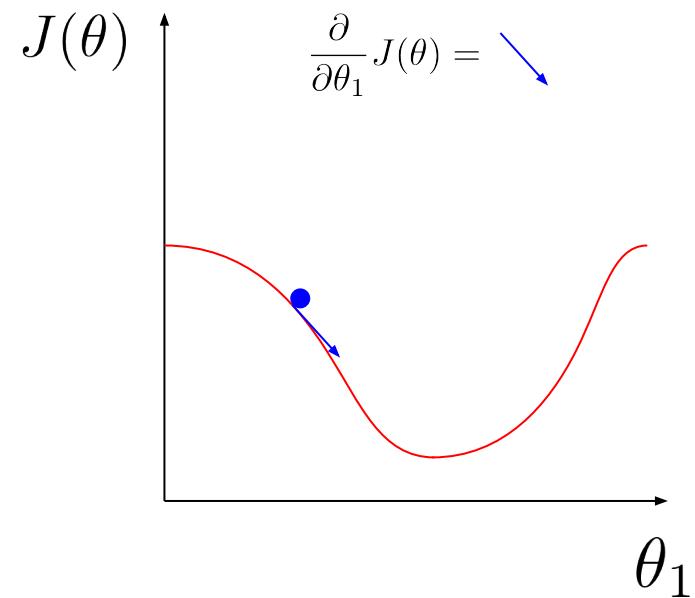
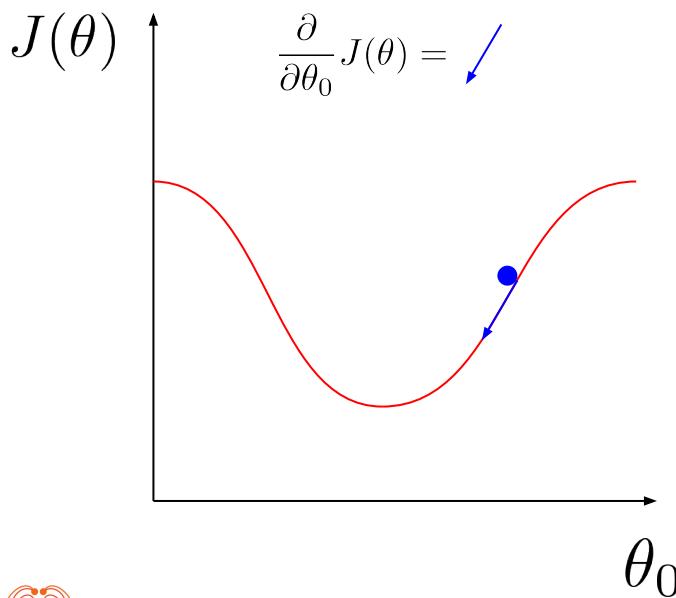
# Gradient descent

We find the direction of steepest slope, then take one step in that direction.



# Gradient descent

Finding the slope for each individual parameter



# Gradient descent

To find the slope, we compute the derivative of the cost (gradient) with respect to a single parameter.

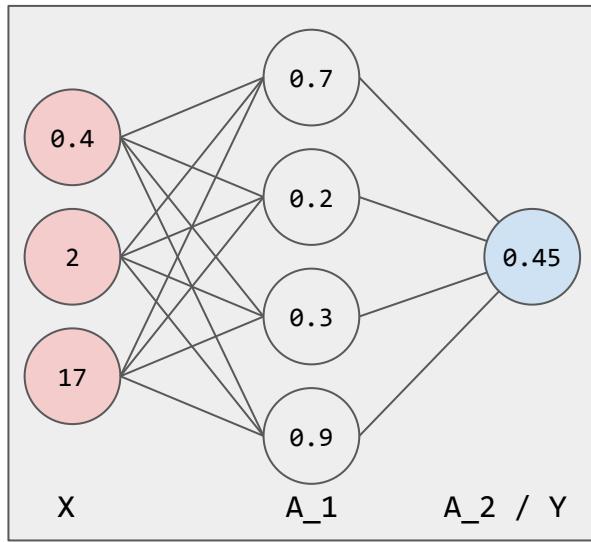
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

Also written  $\nabla J(\theta)$

The diagram illustrates the gradient descent update rule. It shows the formula  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$ . Arrows point from the labels to the corresponding parts of the formula: 'Scalar learning rate' points to  $\alpha$ , 'Individual weights' points to  $\theta_j$ , 'Vector of weights' points to the bracket under  $J(\theta)$ , and 'Cost/objective/loss function' points to  $J(\theta)$ .

# Gradient descent for Neural Networks

Goal: update the weights and biases such that the cost function will output a smaller value  
(i.e. the difference actual and predicted values will be minimized)



Repeat for many iterations (training steps):

Compute forward pass and the cost function,  $J$

$$W^{[1]} := W^{[1]} - \alpha \frac{\partial J}{\partial W^{[1]}}$$

$$b^{[1]} := b^{[1]} - \alpha \frac{\partial J}{\partial b^{[1]}}$$

$$W^{[2]} := W^{[2]} - \alpha \frac{\partial J}{\partial W^{[2]}}$$

$$b^{[2]} := b^{[2]} - \alpha \frac{\partial J}{\partial b^{[2]}}$$

# Chain Rule

1. To get the direction that a parameter must change in order reduce prediction error, we use the **chain rule** from calculus.

$$\frac{d}{dx} f(g(x)) = \overbrace{f'(g(x))}^{\text{outer function}} \underbrace{g'(x)}_{\text{inner function}}$$

$$\frac{\partial J}{\partial W^{[2]}} = \frac{\partial J}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial W^{[2]}}$$

$$\frac{\partial J}{\partial W^{[1]}} = \frac{\partial J}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial W^{[1]}}$$



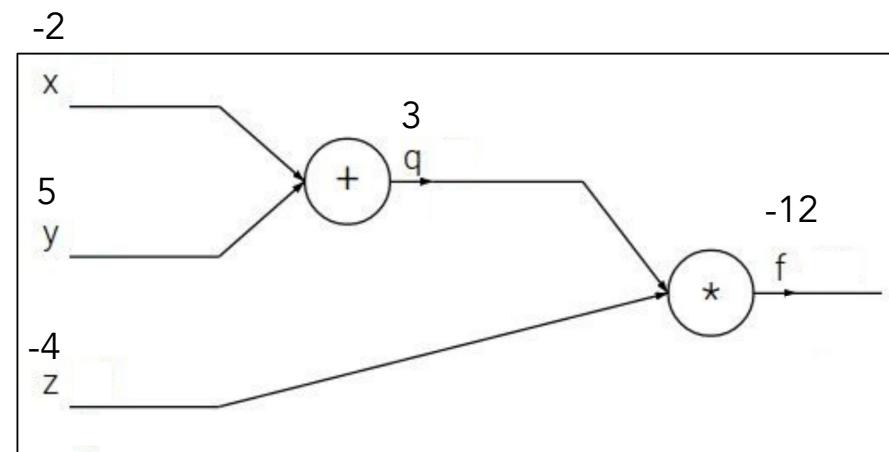
# Chain Rule example

$$f(x, y, z) = (x + y)z$$

Eg.  $x = -2, y = 5, z = -4$

$$q = x + y$$

$$f = qz$$



# Chain Rule example

$$f(x, y, z) = (x + y)z$$

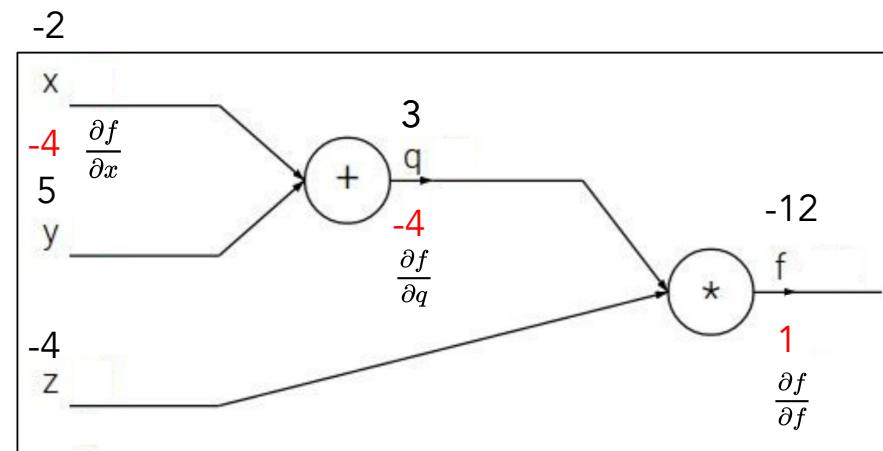
Eg.  $x = -2, y = 5, z = -4$

$$q = x + y$$

$$f = qz$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

$$= z \times 1 = (-4) \times 1 = -4$$



$$J(\hat{y}, y) = \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

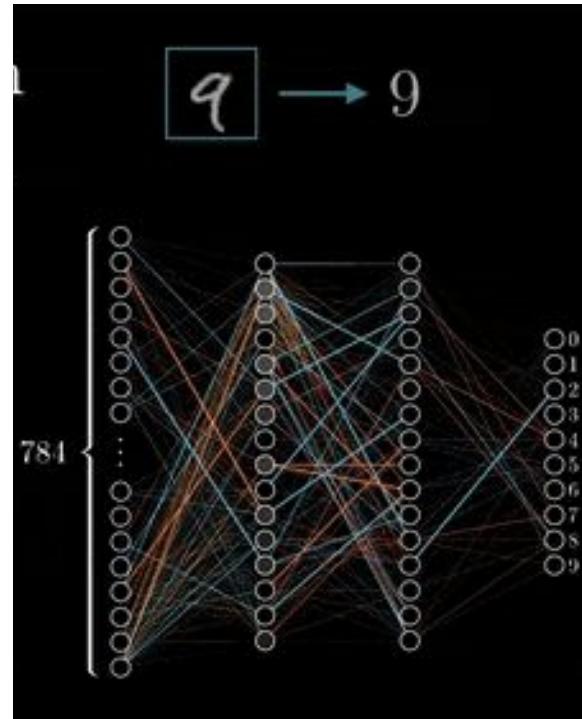
$$L(\hat{y}, y) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

$$a^{[2]} = g^{[2]}(z^{[2]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[1]} = W^{[1]}x + b^{[1]}$$



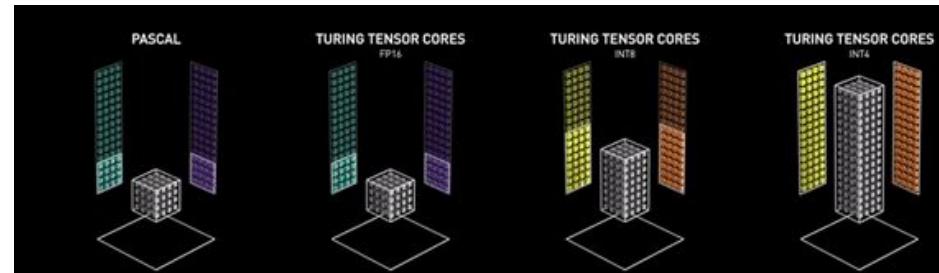


# Speeding up with Linear Algebra

Vectors and matrices

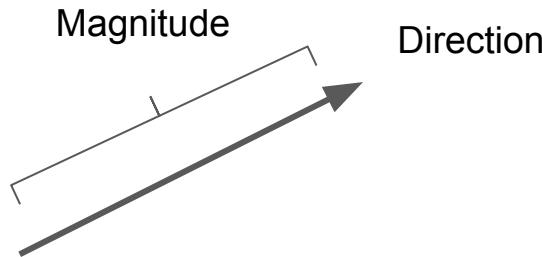
# Vectorization

1. Recent models require massive amounts of computational power
2. Therefore obtaining the most efficiency out of your model is crucial
3. Linear algebra libraries are becoming more efficient, and neural networks can benefit from this outcome
  - Matrix multiplication in Numpy
  - Tensor operations in TensorFlow



# Vectors and Matrices

## Vector



Feature vector

$$x = \begin{bmatrix} 0.4 \\ 2 \\ 17 \end{bmatrix}$$

## Matrix

Weight Matrix

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}$$

rows

columns

# Vectors and Matrices

Inner product (dot product)

$$v \cdot w = v^T w = [1 \ 2 \ 3] \begin{bmatrix} 4 \\ 1 \\ 3 \end{bmatrix} = 1 \times 4 + 2 \times 1 + 3 \times 3 = 15$$

Matrix Multiplication

Matrix vector multiplication

$$\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} ax + by + cz \\ dx + ey + fz \\ gx + hy + iz \end{bmatrix}$$

3x3      3x1      =      3x1

$$\begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 0 \\ 2 & 3 & 4 \end{bmatrix} \times \begin{bmatrix} 2 & 5 \\ 6 & 7 \\ 1 & 8 \end{bmatrix}$$

-    +                  -    +

-    +                  -    +

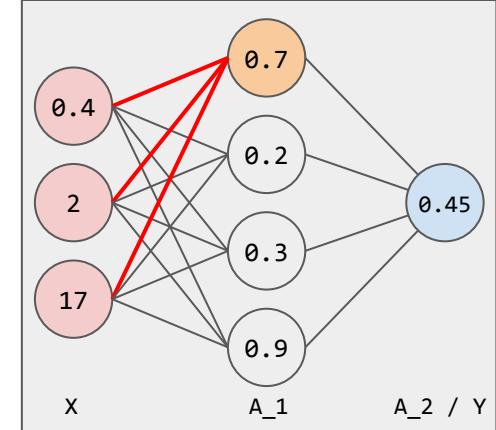
► Multiply

## Vectorizing within one activation

$$z^{[l]} = \left( w_1^{[l]} a_1^{[l-1]} + \dots + w_n^{[l]} a_n^{[l-1]} + b^{[l]} \right)$$

$$= \begin{bmatrix} w_1^{[l]} \\ \vdots \\ w_n^{[l]} \end{bmatrix}^T \begin{bmatrix} a_1^{[l-1]} \\ \vdots \\ a_n^{[l-1]} \end{bmatrix} + b^{[l]} = w^{[l]T} a^{[l-1]} + b^{[l]}$$

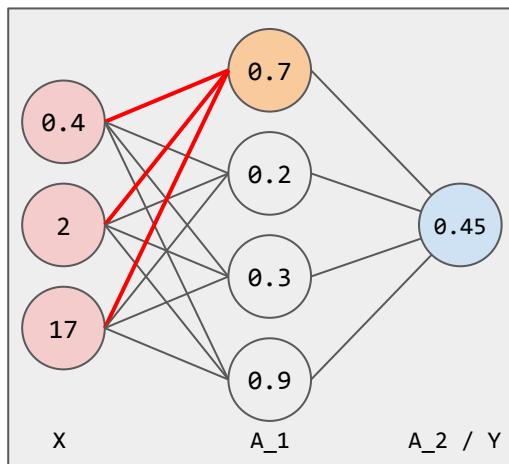
$$a^{[l]} = g^{[l]}(z^{[l]}) = g^{[l]}(w^{[l]T} a^{[l-1]} + b^{[l]})$$



# Vectorizing within one activation (example)

$$z^{[l]} = \left( w_1^{[l]} a_1^{[l-1]} + \dots + w_n^{[l]} a_n^{[l-1]} + b^{[l]} \right) = w^{[l]T} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]}) = g^{[l]}(w^{[l]T} a^{[l-1]} + b^{[l]})$$

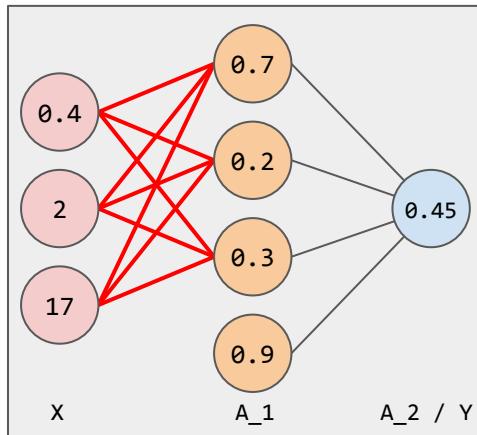


$$\begin{aligned} z_1^{[1]} &= w_1^{[1]} x_1 + w_2^{[1]} x_2 + w_3^{[1]} x_3 + b \\ &= \underbrace{0.5 * 0.4 + 0.1 \times 2 + 0.0058 \times 17}_{\text{sum}} + 0.2 = 0.7 \end{aligned}$$

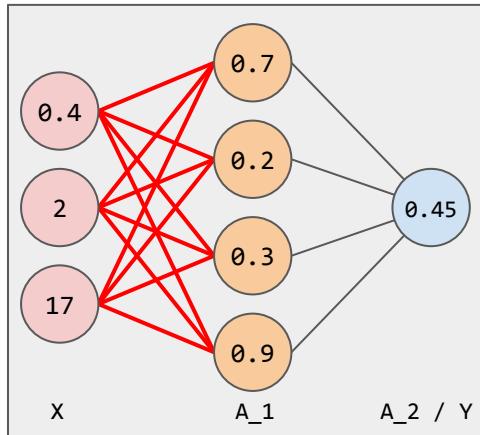
$$\begin{bmatrix} 0.5 \\ 0.1 \\ 0.0058 \end{bmatrix}^T \begin{bmatrix} 0.4 \\ 2 \\ 17 \end{bmatrix} + 0.2 = 0.7$$

$$a_1^{[1]} = g(0.7) = \text{ReLU}(0.7) = 0.7$$

# Vectorizing across all activations

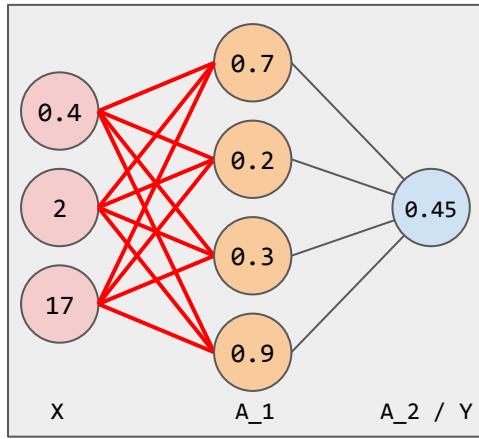


# Vectorizing across all activations



$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, & a_1^{[1]} &= g^{[1]}(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, & a_2^{[1]} &= g^{[1]}(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, & a_3^{[1]} &= g^{[1]}(z_3^{[1]}) \\ z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, & a_4^{[1]} &= g^{[1]}(z_4^{[1]}) \end{aligned}$$

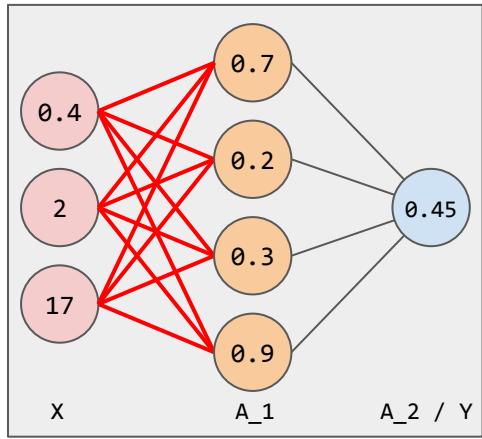
# Vectorizing across all activations



$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} \boxed{x} + b_1^{[1]}, & a_1^{[1]} &= g^{[1]}(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} \boxed{x} + b_2^{[1]}, & a_2^{[1]} &= g^{[1]}(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} \boxed{x} + b_3^{[1]}, & a_3^{[1]} &= g^{[1]}(z_3^{[1]}) \\ z_4^{[1]} &= w_4^{[1]T} \boxed{x} + b_4^{[1]}, & a_4^{[1]} &= g^{[1]}(z_4^{[1]}) \end{aligned}$$

$$z^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]T}x + b_1^{[1]} \\ w_2^{[1]T}x + b_2^{[1]} \\ w_3^{[1]T}x + b_3^{[1]} \\ w_4^{[1]T}x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

# Vectorizing across all activations



$$z^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]T}x + b_1^{[1]} \\ w_2^{[1]T}x + b_2^{[1]} \\ w_3^{[1]T}x + b_3^{[1]} \\ w_4^{[1]T}x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$$z^{[l]}(i) = W^{[l]}a^{[l-1]}(i) + b^{[l]}, \quad a^{[l]}(i) = g^{[l]}(z^{[l]}(i))$$

for i = 1 to m:

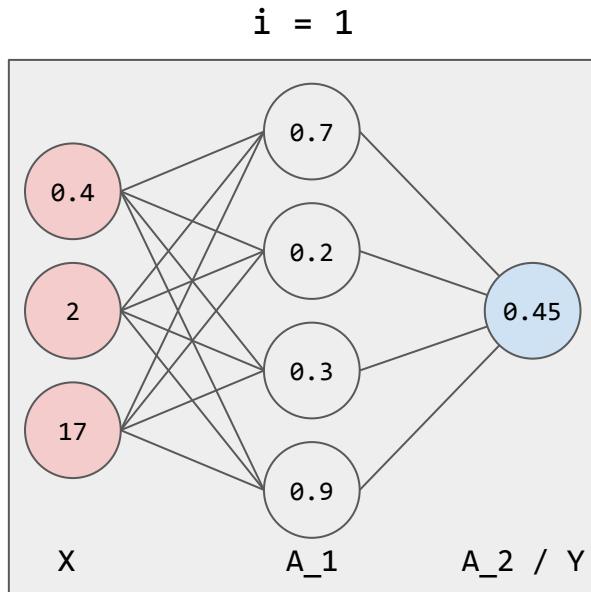
$$z^{[1]}(i) = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1]}(i) = \sigma(z^{[1]}(i))$$

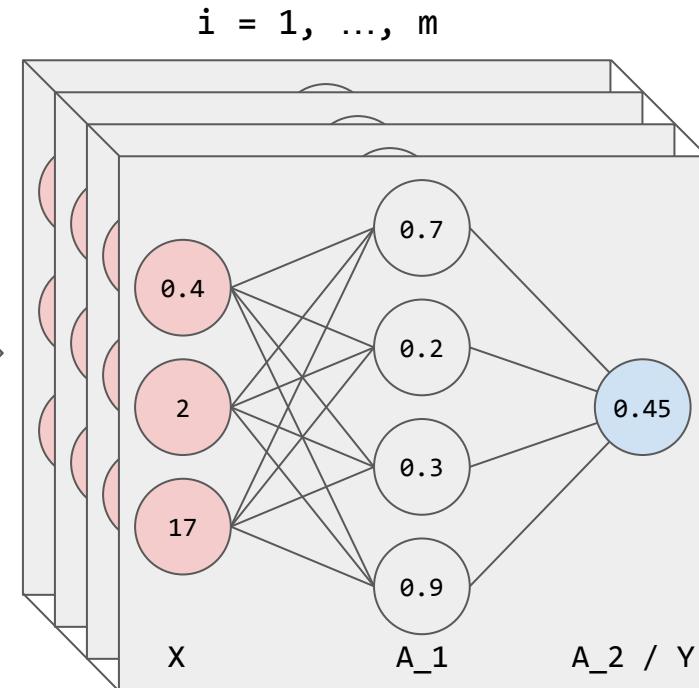
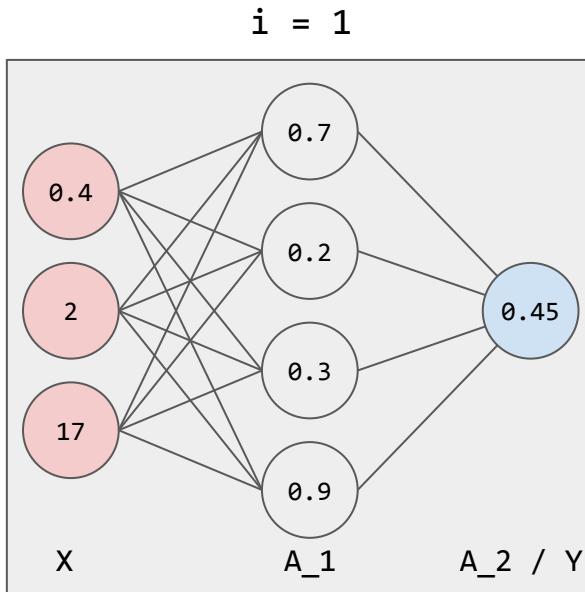
$$z^{[2]}(i) = W^{[2]}a^{[1]}(i) + b^{[2]}$$

$$a^{[2]}(i) = \sigma(z^{[2]}(i))$$

# Single Example Neural Network

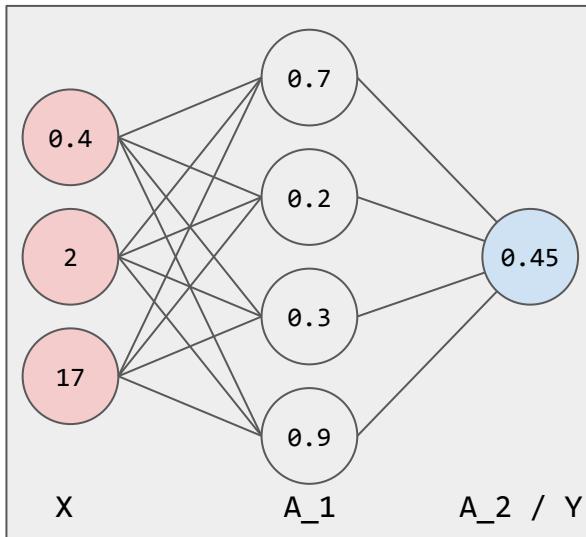


# Vectorizing across multiple examples



# Vectorizing across multiple examples

**Algorithm: forward propagate through each example**



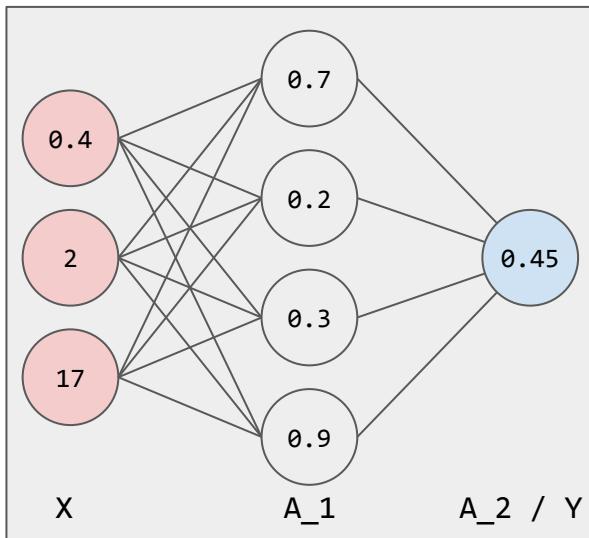
```
for i = 1 to m:  
     $z^{[1]}(i) = W^{[1]}x^{(i)} + b^{[1]}$   
     $a^{[1]}(i) = \sigma(z^{[1]}(i))$   
     $z^{[2]}(i) = W^{[2]}a^{[1]}(i) + b^{[2]}$   
     $a^{[2]}(i) = \sigma(z^{[2]}(i))$   

$$J(\hat{y}, y) = \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

```

# Vectorizing across multiple examples

**Algorithm: forward propagate through each example**



$$z^{[l]}(i) = W^{[l]}a^{[l-1](i)} + b^{[l]}, \quad a^{[l]}(i) = g^{[l]}(z^{[l]}(i))$$

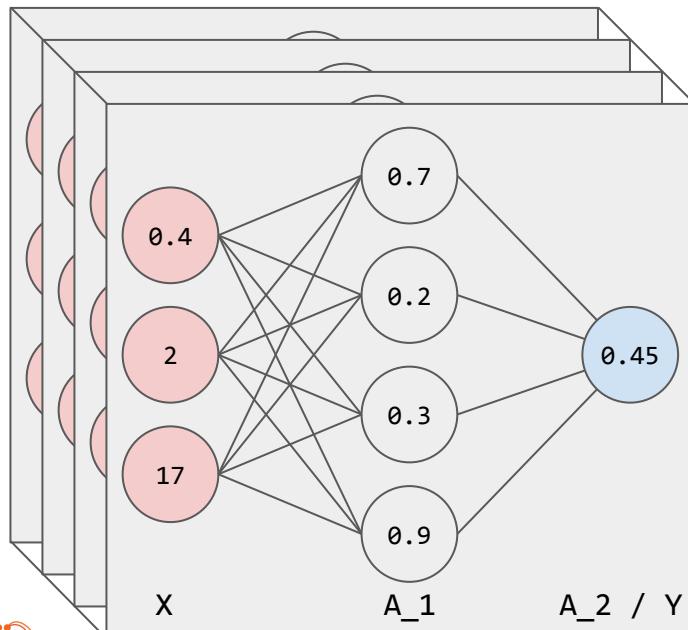
Diagram illustrating the forward propagation equation:

$$z^{[l]} = \begin{matrix} \text{---} \\ | \\ \text{---} \\ | \\ \text{---} \\ | \\ \text{---} \end{matrix} = \begin{matrix} n_h^{[l-1]} \\ \text{---} \\ | \\ \text{---} \\ | \\ \text{---} \\ | \\ \text{---} \end{matrix} \times \begin{matrix} \text{---} \\ | \\ \text{---} \\ | \\ \text{---} \\ | \\ \text{---} \end{matrix} + \begin{matrix} n_h^{[l]} \\ \text{---} \\ | \\ \text{---} \\ | \\ \text{---} \\ | \\ \text{---} \end{matrix}$$

Where  $z^{[l]}$  is the vector of pre-activations,  $W^{[l]}$  is the weight matrix,  $a^{[l-1]}$  is the previous activation vector, and  $b^{[l]}$  is the bias vector.

# Vectorizing across multiple examples

**Vectorized algorithm: forward propagate through multiple examples in parallel**

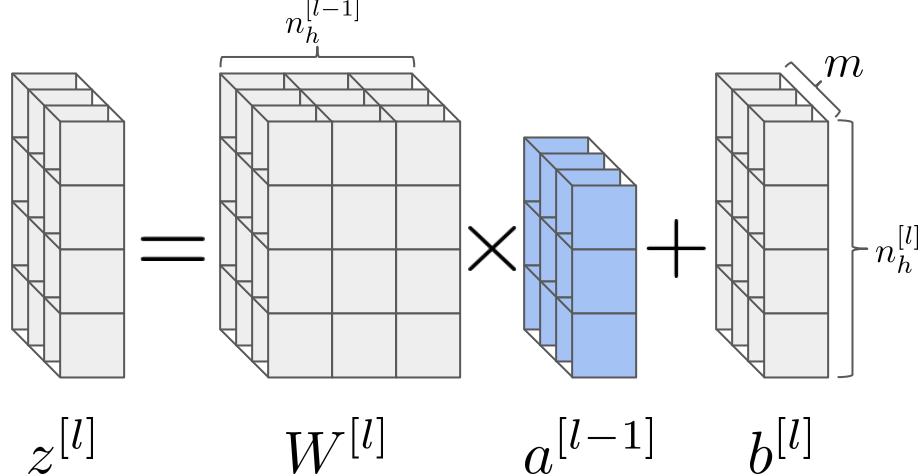


$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$

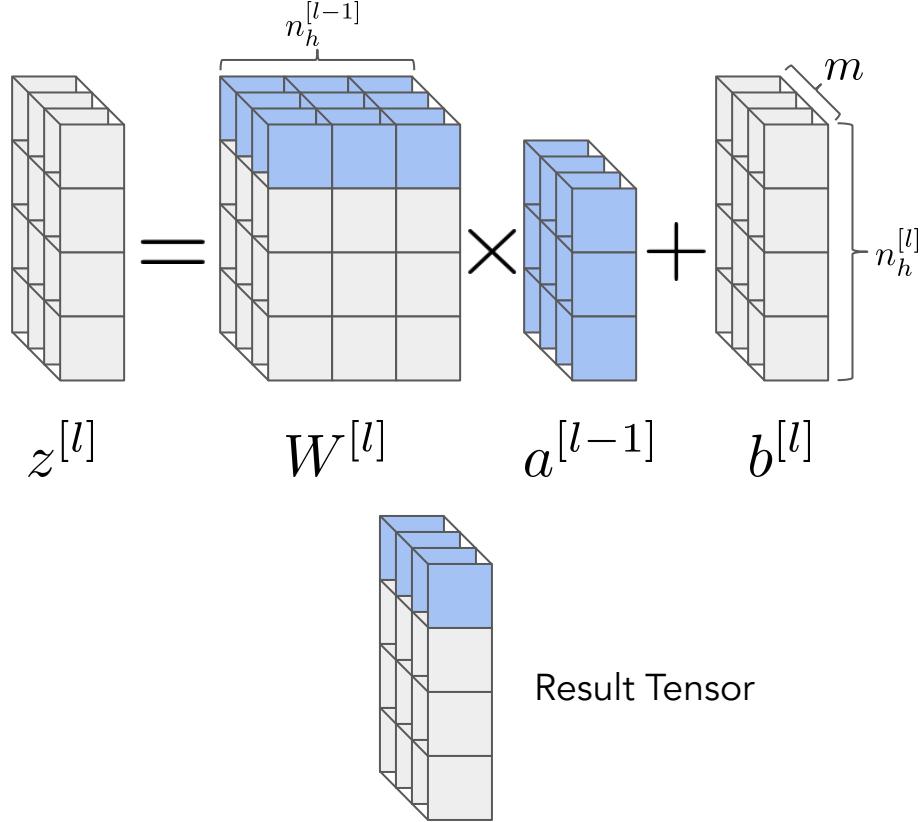
$$z^{[l]} = \begin{matrix} n_h^{[l-1]} \\ \times \\ W^{[l]} \\ + \\ b^{[l]} \end{matrix} = \begin{matrix} m \\ \times \\ a^{[l-1]} \\ + \\ b^{[l]} \end{matrix}$$

The diagram shows the mathematical representation of the forward pass. It consists of three main components: a weight matrix  $W^{[l]}$  (represented as a 3D cube), a bias vector  $b^{[l]}$  (represented as a 2D column), and an activation function  $g^{[l]}$ . The input  $a^{[l-1]}$  is multiplied by  $W^{[l]}$  and then added to  $b^{[l]}$  to produce the output  $z^{[l]}$ . The dimensions of the weight matrix  $W^{[l]}$  are indicated as  $n_h^{[l-1]} \times m$ , where  $n_h^{[l-1]}$  is the number of nodes in the previous layer and  $m$  is the number of nodes in the current layer.

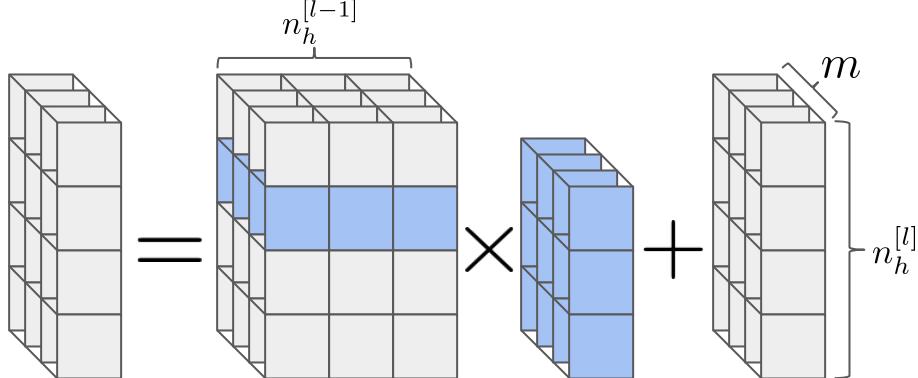
$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$



$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$



$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$

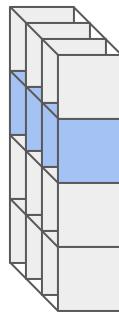


$z^{[l]}$

$W^{[l]}$

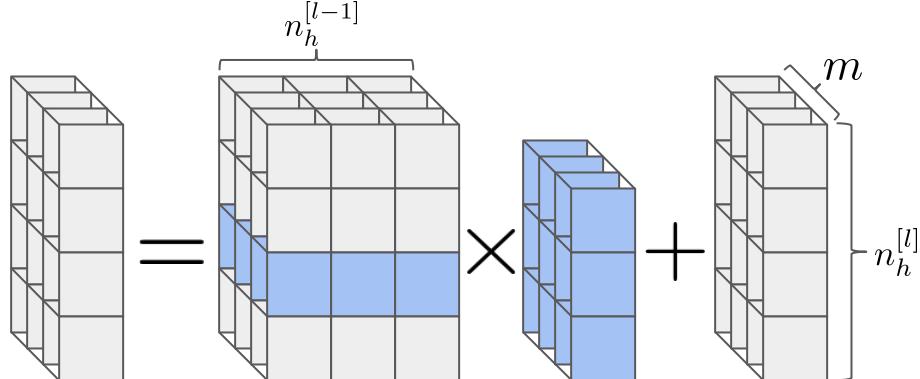
$a^{[l-1]}$

$b^{[l]}$

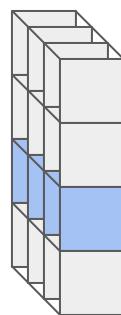


Result Tensor

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$

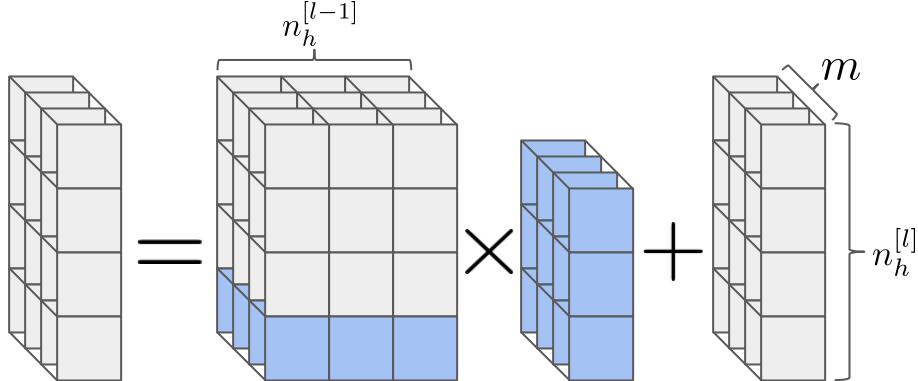
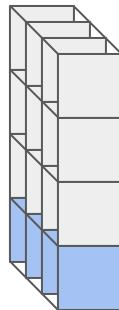


$$z^{[l]} \quad W^{[l]} \quad a^{[l-1]} \quad b^{[l]}$$



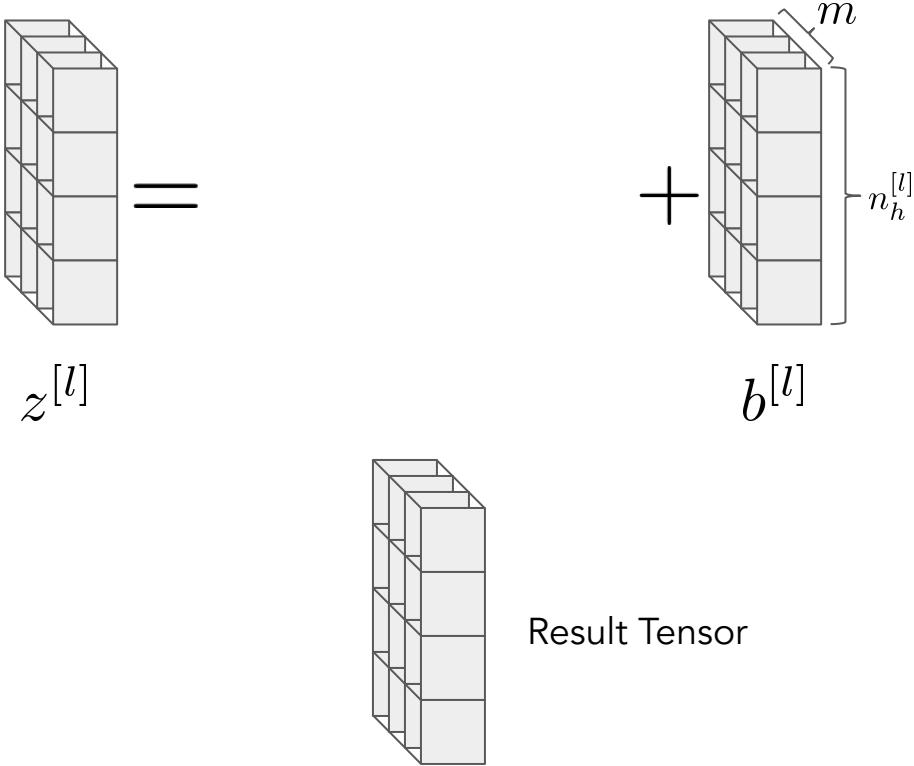
Result Tensor

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$

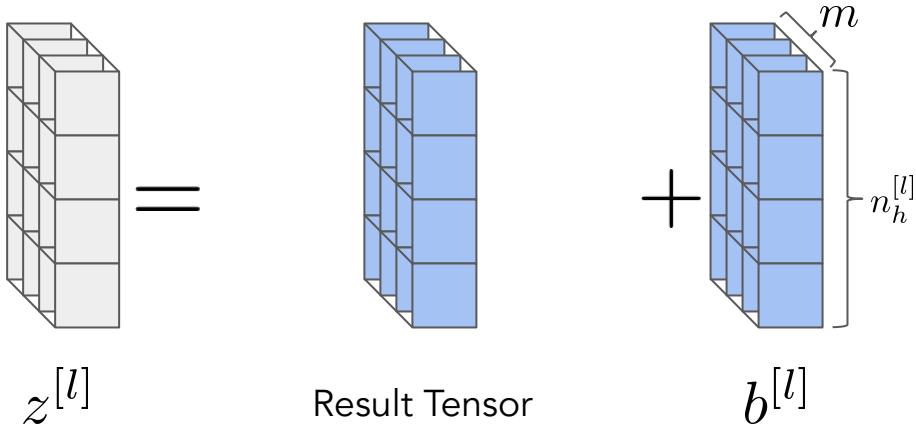
 $z^{[l]}$  $W^{[l]}$  $a^{[l-1]}$  $b^{[l]}$ 

Result Tensor

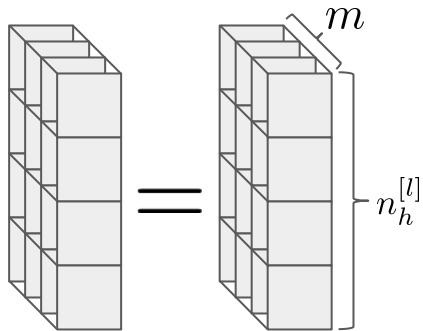
$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$



$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$



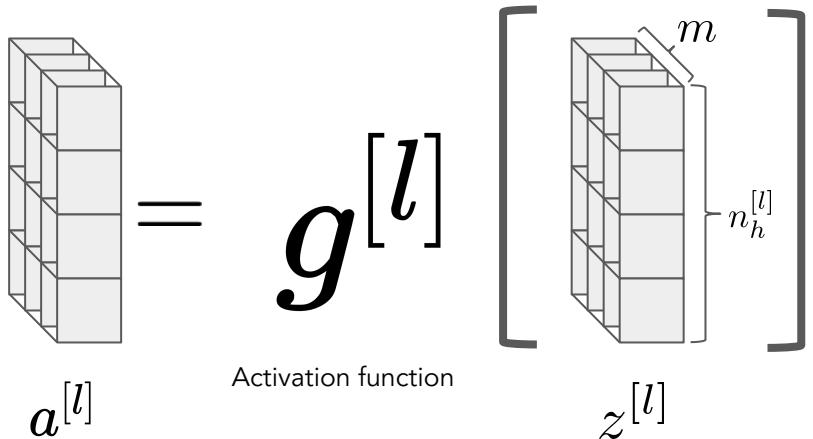
$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$



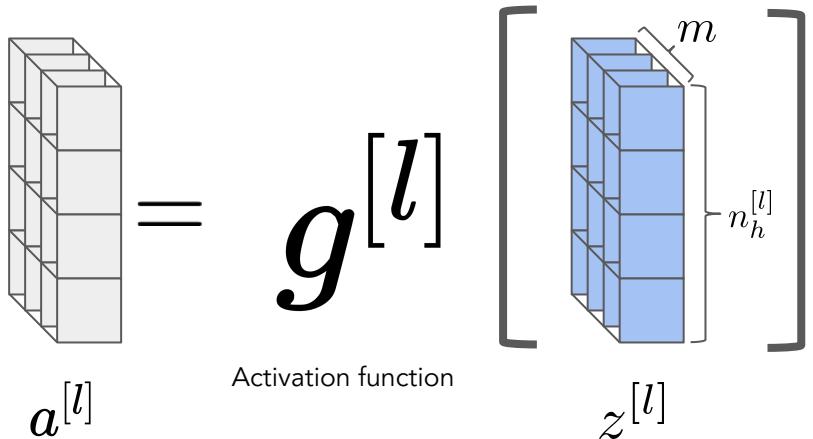
$z^{[l]}$       Result Tensor



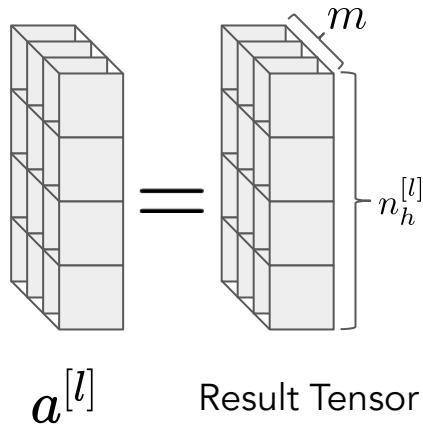
$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$



$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$



$$z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$



# Vectorizing Gradient Descent

**Going from:** Iteratively changing each weight by the derivative of the cost function

**To:** In parallel, simultaneously changing each weight by the derivative of the cost function



# Vectorizing Gradient Descent

**Going from:** Iteratively changing each weight by the derivative of the cost function

**To:** In parallel, simultaneously changing each weight by the derivative of the cost function

**Previously,**  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$

**Now, W is the weight matrix for layer l**

$$W^{[l]} := W^{[l]} - \alpha \frac{\partial}{\partial W^{[l]}} J(\theta)$$

Note similarly we do this for the bias vector as well

# Recap: Gradient descent for Neural Networks

Repeat for many iterations (training steps):

Compute forward pass and the cost function,  $\mathbf{J}$

$$W^{[1]} := W^{[1]} - \alpha \frac{\partial J}{\partial W^{[1]}}$$

$$b^{[1]} := b^{[1]} - \alpha \frac{\partial J}{\partial b^{[1]}}$$

$$W^{[2]} := W^{[2]} - \alpha \frac{\partial J}{\partial W^{[2]}}$$

$$b^{[2]} := b^{[2]} - \alpha \frac{\partial J}{\partial b^{[2]}}$$



# Derivatives for a simple Neural Network

$$dZ^{[2]} = A^{[2]} - Y$$

Note:  $dZ^{[1]} \rightarrow \frac{\partial J}{\partial Z^{[1]}}$

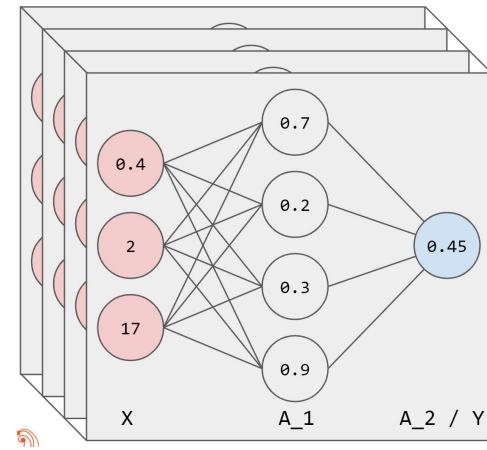
$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis = 1, keepdims = True)$$

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

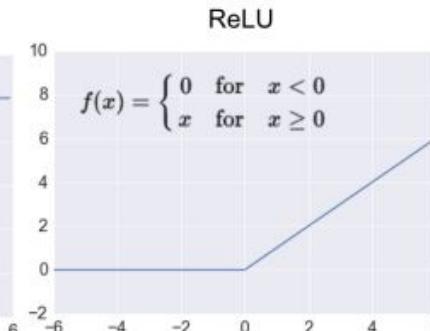
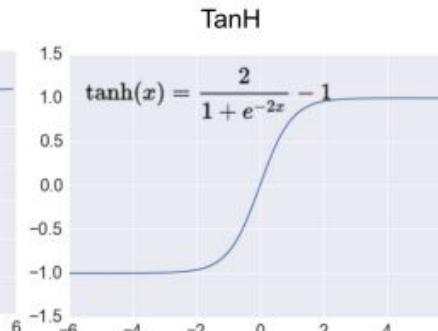
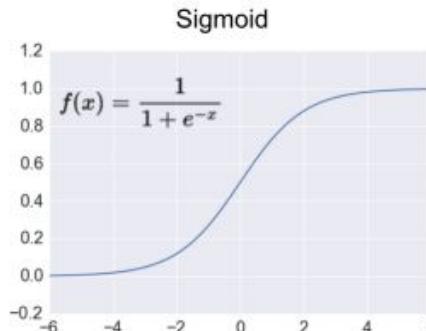
$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$$



# Derivatives for activation functions

$$g'(x) = \begin{cases} \frac{d}{dx} \sigma(x) = \sigma(x)(1 - \sigma(x)) \\ \frac{d}{dx} \tanh(x) = 1 - \tanh^2(x) \\ \frac{d}{dx} \text{ReLU}(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x > 0 \end{cases} \end{cases}$$





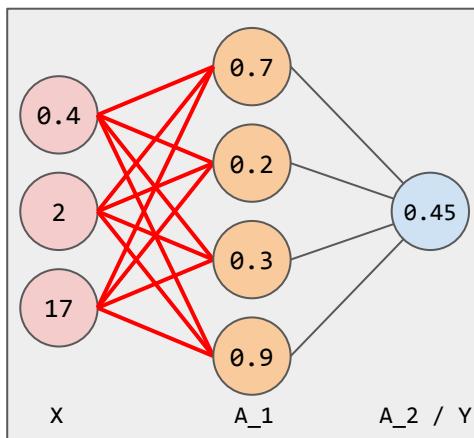
# Deep Neural Networks

The general case

# Abstracting Layers

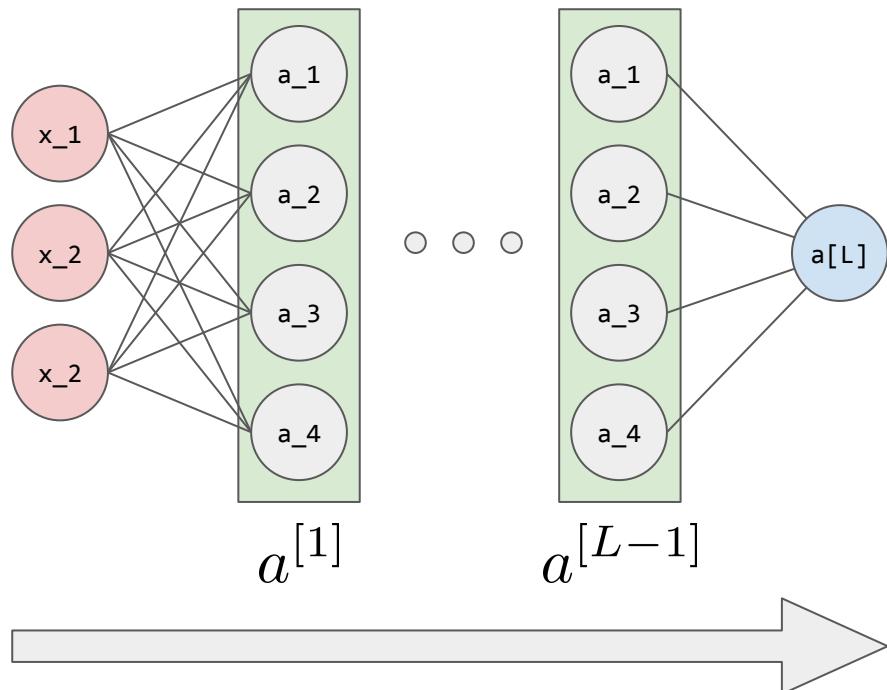
A layer is an abstract generalization of a transformation function on the previous activations.

$$a^{[l]} = f^{[l]}(a^{[l-1]})$$



$$a^{[1]} = f^{[1]}(\mathbf{x})$$

# Forward Propagation

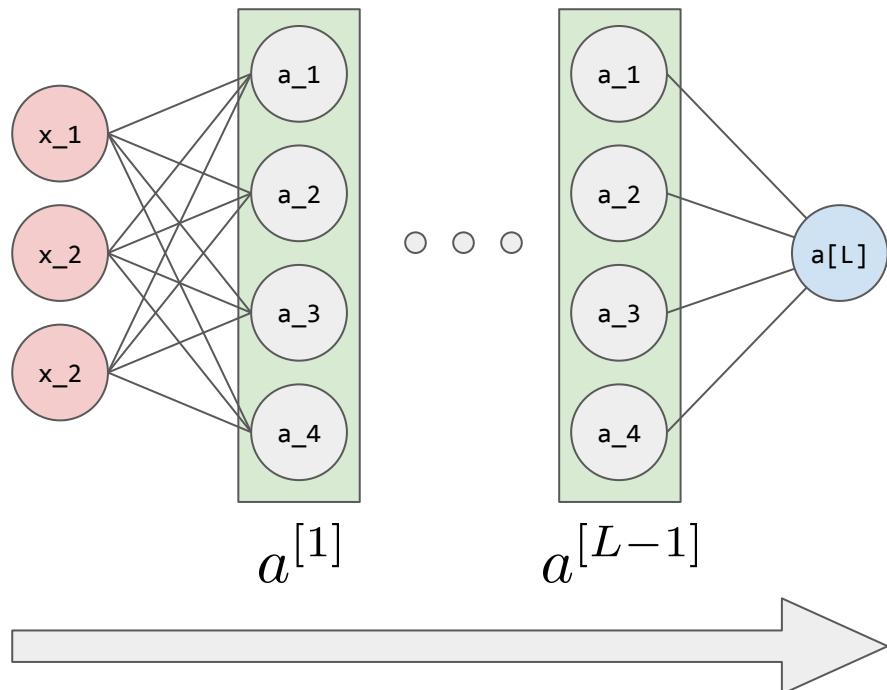


Repeat L-1 times

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$

$= \begin{matrix} z^{[l]} \\ \times \\ W^{[l]} \\ + \\ b^{[l]} \end{matrix} = \begin{matrix} n_h^{[l-1]} \\ \times \\ m \\ + \\ n_h^{[l]} \end{matrix}$

# Layered Forward Propagation



Repeat  $L-1$  times

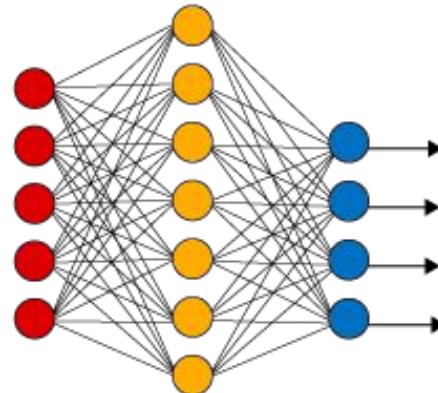
$$a^{[l]} = f^{[l]}(a^{[l-1]})$$

$= f^{[l]} \left [ \begin{array}{c} \text{[matrix icon]} \\ \vdots \\ \text{[matrix icon]} \end{array} \right ]$

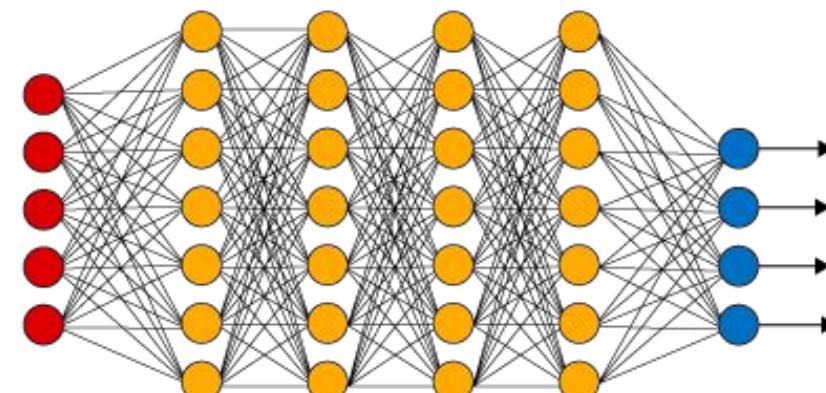
# Deep Neural Networks

1. Just a neural network with more layers
2. The number of layers is denoted as  $L$

Simple Neural Network



Deep Learning Neural Network



● Input Layer

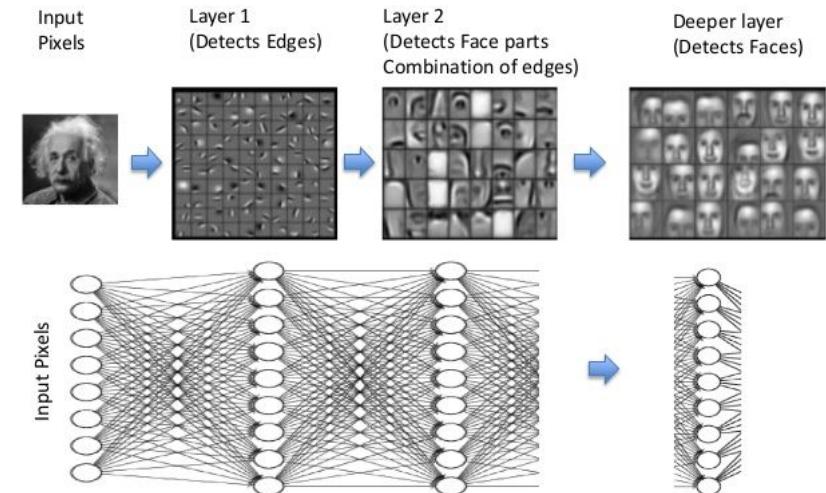
● Hidden Layer

● Output Layer

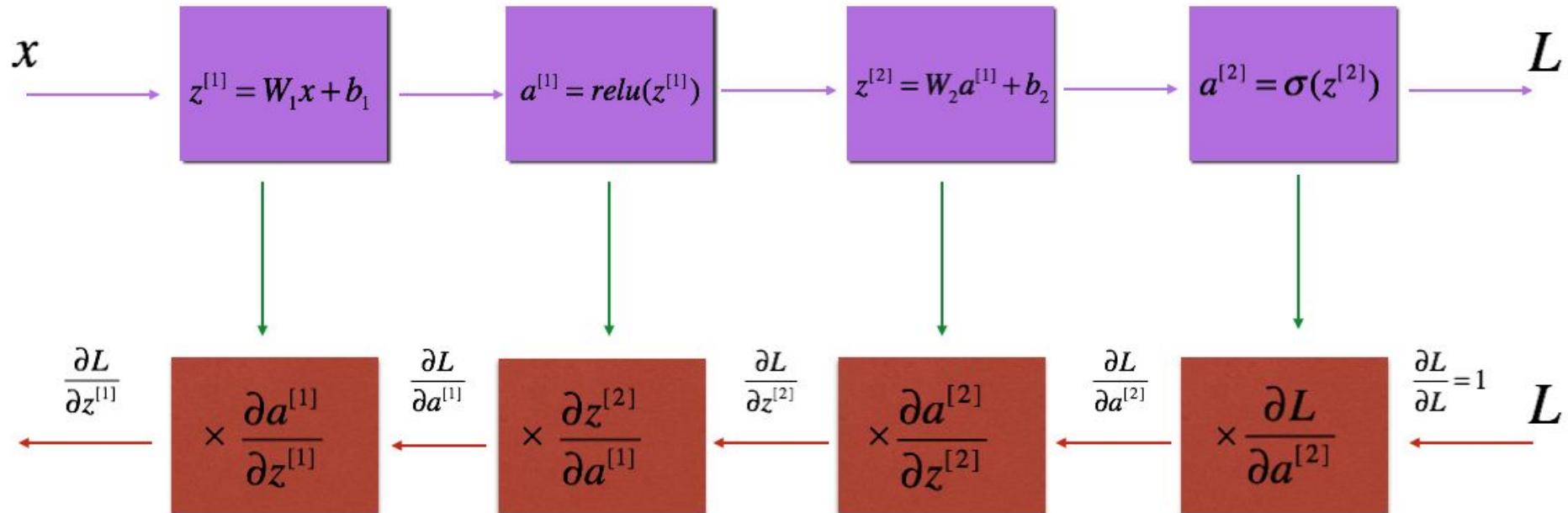
# Why do we need more layers?

1. Lower layers learn more basic features of the input (edge detection)
2. Deeper layers build off of previously learned features to conceptualize more complex features (face detection)
3. Adding more layers allows the network to approximate more complex functions

Feature Learning/Representation Learning  
(Ex. Face Detection)



# Backward Propagation



# Backward Propagation

$$dZ^{[L]} = A^{[L]} - Y$$

$$dW^{[L]} = \frac{1}{m} dZ^{[L]} A^{[L-1]T}$$

$$db^{[L]} = \frac{1}{m} \text{np.sum}(dZ^{[L]}, \text{axis} = 1, \text{keepdims} = \text{True})$$

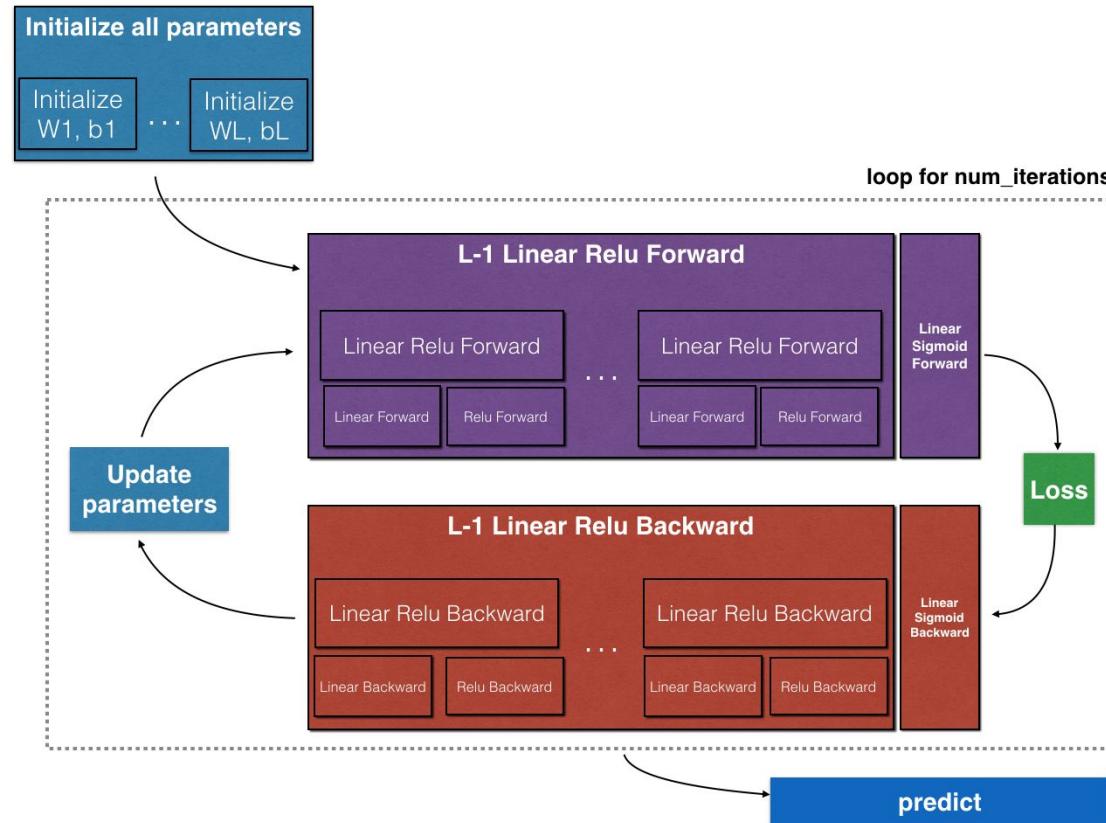
$$dZ^{[l]} = W^{[l+1]T} dZ^{[l+1]} * g^{[l]'}(Z^{[l]})$$

for l = L-1 to 1

$$dW^{[l]} = \frac{1}{m} dZ^{[l]} A^{[l-1]T}$$

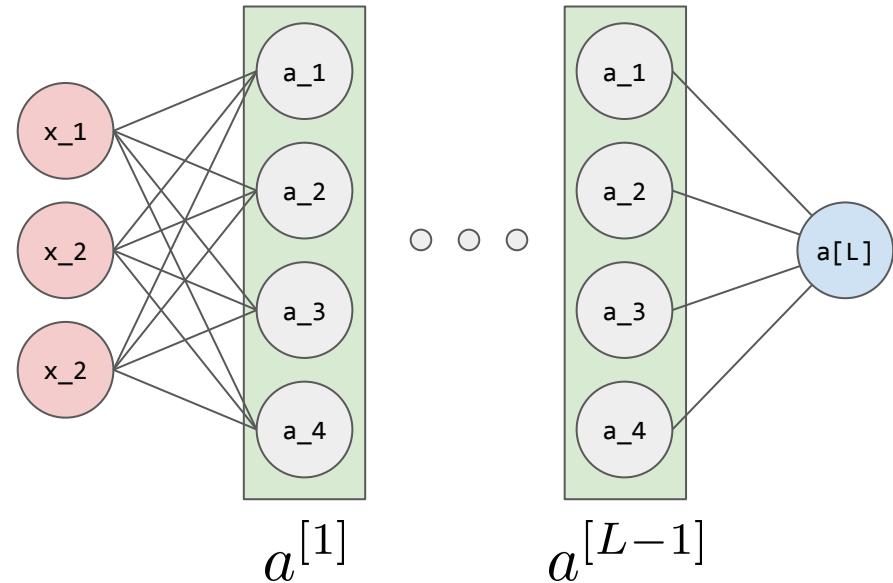
$$db^{[l]} = \frac{1}{m} \text{np.sum}(dZ^{[l]}, \text{axis} = 1, \text{keepdims} = \text{True})$$

# Building blocks of deep neural networks



# Summary

1. Forward propagate
2. Compute the cost
3. Back propagate
4. Update weights



# References and Further Reading

1. Deep Learning Specialization by Andrew Ng: <https://wwwdeeplearning.ai/>
2. Open AI five: <https://blog.openai.com/openai-five/>
3. AlphaGo Zero: <https://deepmind.com/blog/alphago-zero-learning-scratch/>
4. Google Colab, Free Cloud GPUs: <https://colab.research.google.com>
5. BU MIC Youtube Channel: <https://www.youtube.com/channel/UCWZLArtz3qv7bVaughcl1bA>
6. Feedback form for this workshop: <https://goo.gl/forms/i9dNSEzmGjFKbsXD3>
7. Deep Learning Applications: [Forbes Deep Learning Applications](#)
8. Deep Learning vs. ML: <https://www.zendesk.com/blog/machine-learning-and-deep-learning/>
9. AlexNet: <https://en.wikipedia.org/wiki/AlexNet>

