

# Neural Networks & Deep Learning

## What to learn?

NN, DL

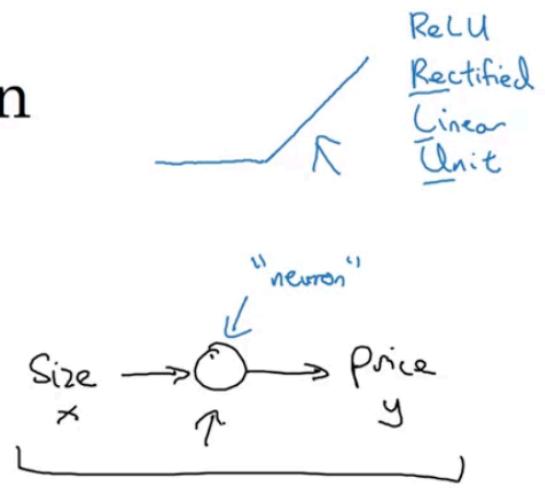
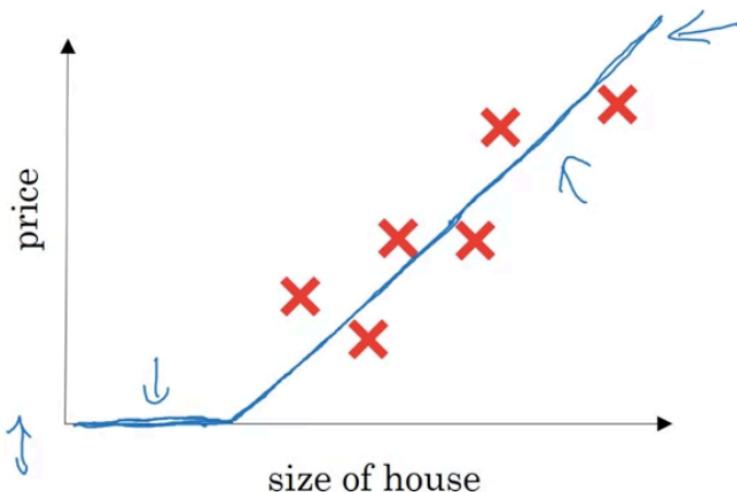
Improvements: Hyperparameters, Regularization, Optimization

Structuring projects

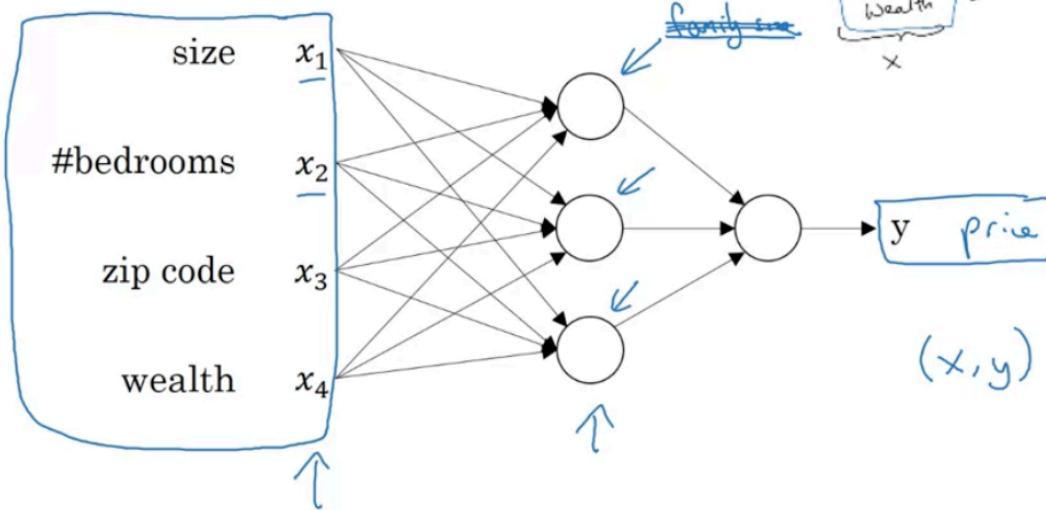
Convolutional NNs, CNNs

Sequence Models, NLP, RNN, LSTM

## Housing Price Prediction



## Housing Price Prediction



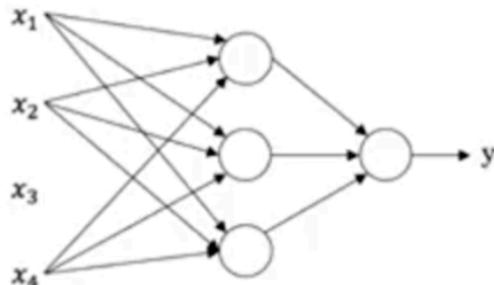
Relatively standard for real estate: Standard NN

Image: CNN

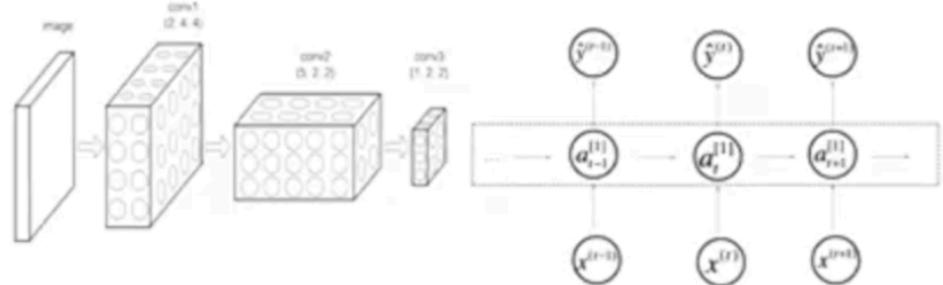
Sequence Data: (Audio) One dim. type of sequence data: RNN

Image + Radar: Custom/Hybrid NNs

# Neural Network examples



Standard NN



Convolutional NN

Recurrent NN

Structured Data: Databases -> clear labels, size, bedroom number, price, age, id...

Unstructured Data: Audio, Images, Text

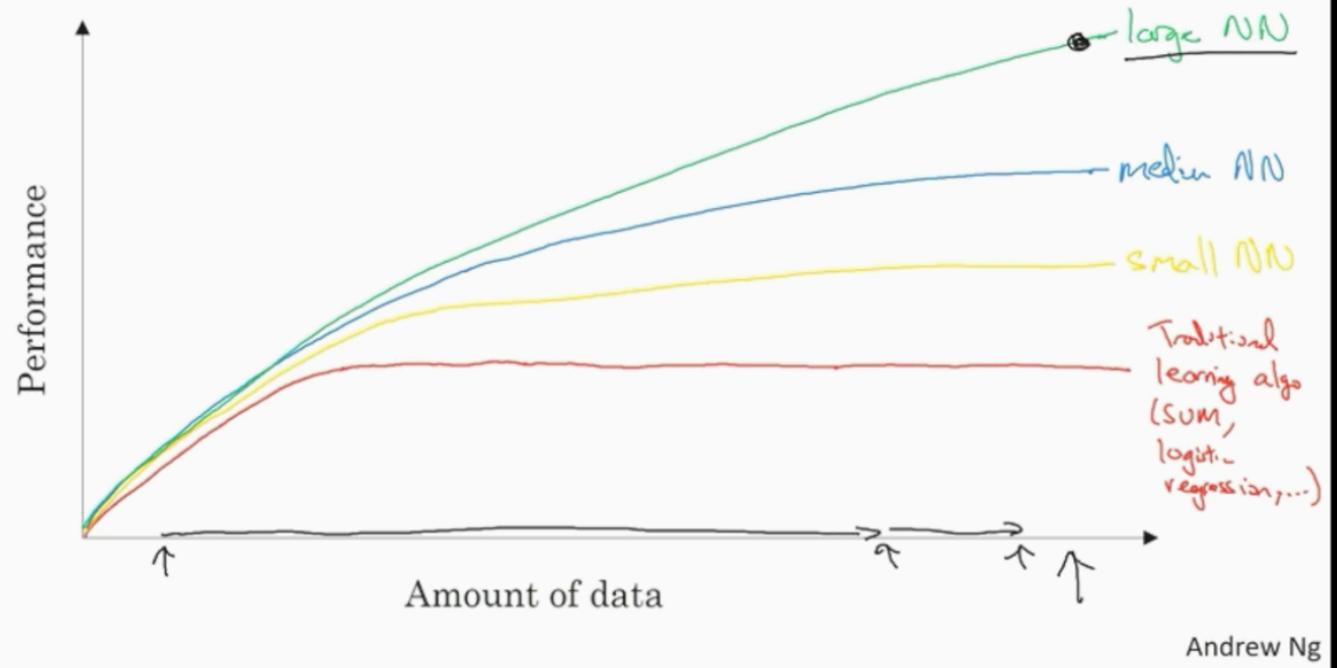
## Why now?

Scale drives DL

Traditional learning algo: Support Vector Machine, Logistic Regression

- learning curve flattens (plateaus after some amount of data)  
We continue collecting data

## Scale drives deep learning progress



just improving scale on both ends helps (mostly until the hitting a limit)

for small data, the performance of model architectures are not clearly listed. Depends on the specific situation

Data - Computation - Algorithm (mostly, enabled the NNs run faster)

Sigmoid fnc (parameters move too slow) -> ReLU (rectified last unit)  
makes the gradient descent much faster

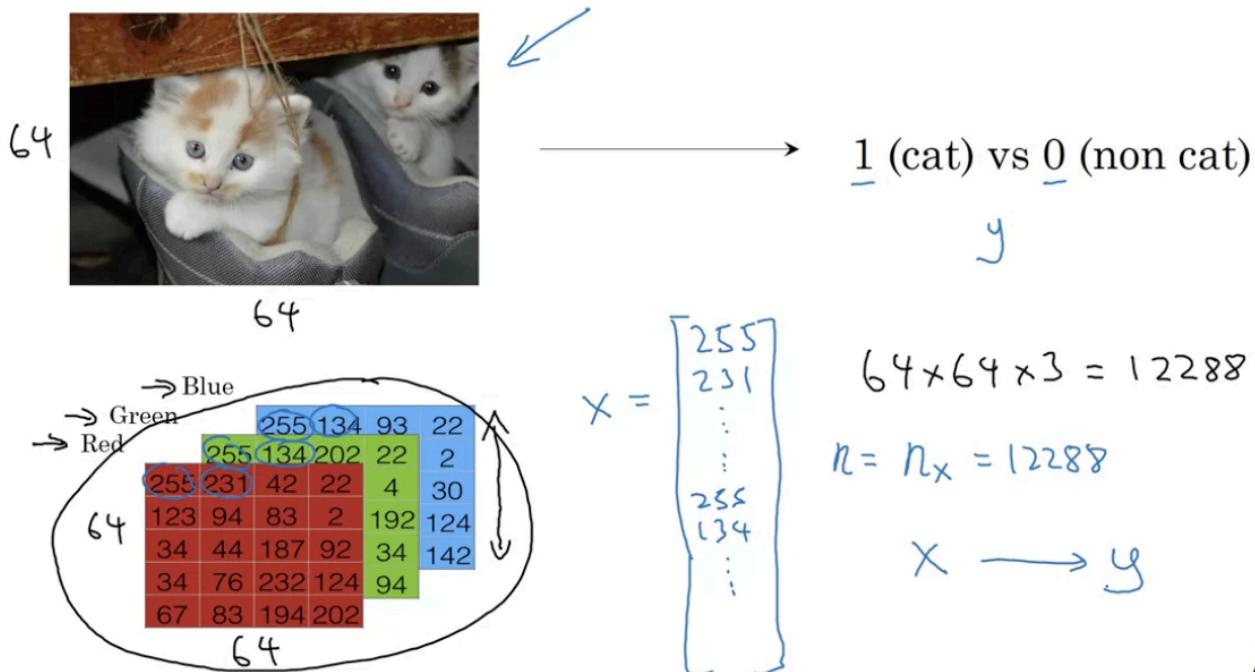
- Iterative process: Idea -> Code -> *Experiment* -> Repeat

10 min -> Day -> Month

**Geoffrey Hinton:** Godfather of DL

## Binary Classification

# Binary Classification



## Notation

$$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

$m$  training examples :  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$M = M_{\text{train}}$        $M_{\text{test}} = \# \text{test examples.}$

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}_{n_x \times m}$$

$X \in \mathbb{R}^{n_x \times m}$        $X.\text{shape} = (n_x, m)$

$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$

$Y \in \mathbb{R}^{1 \times m}$

$Y.\text{shape} = (1, m)$

Andrew Ng

## Logistic Regression

given  $x$  - want  $y'$

$x$  picture, chance of cat

linear regression is not a good choice

sigmoid fnc

# Logistic Regression cost function

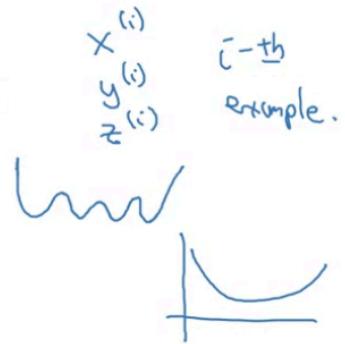
$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \text{ where } \sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}} \quad z^{(i)} = w^T x^{(i)} + b$$

Given  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ , want  $\hat{y}^{(i)} \approx y^{(i)}$ .

Loss (error) function:

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

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



If  $y=1$ :  $L(\hat{y}, y) = -\log \hat{y} \leftarrow$  Want  $\log \hat{y}$  large, want  $\hat{y}$  large.

If  $y=0$ :  $L(\hat{y}, y) = -\log (1-\hat{y}) \leftarrow$  Want  $\log (1-\hat{y})$  large ... want  $\hat{y}$  small

$$\text{Cost function: } J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

Andrew Ng

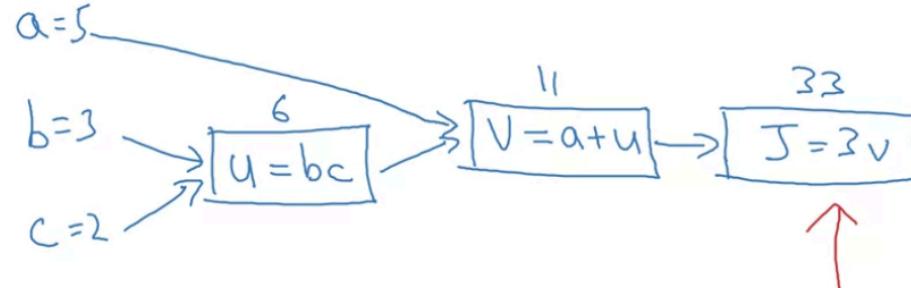
Check: Gradient Descent, Derivatives, etc.

## Computation Graph

$$J(a, b, c) = 3(a + bc) = 3(5 + 3 \cdot 2) = 33$$

$$\begin{array}{c} u = bc \\ v = au \\ J = 3v \end{array}$$

$$\begin{aligned} u &= bc \\ v &= au \\ J &= 3v \end{aligned}$$



Andrew Ng

Logistic Regression derivatives

Gradient Descent on m Examples

## Vectorization

for DL, vectorization is key

even in a small demonstration, for loop is 300x slower compared to vectorized operations.

SIMD: Single instruction, multiple data. Both works for CPU & GPU

whenever possible, avoid explicit for loops.

# Vectors and matrix valued functions

Say you need to apply the exponential operation on every element of a matrix/vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

$\rightarrow u = np.zeros((n, 1))$

~~$\rightarrow \text{for } i \text{ in range}(n):$~~

$\rightarrow u[i] = \text{math.exp}(v[i])$

import numpy as np  
u = np.exp(v) ←  
np.log(u)  
np.abs(u)  
np.maximum(v, 0)  
 $v^{**2}$  ←  $v/v$

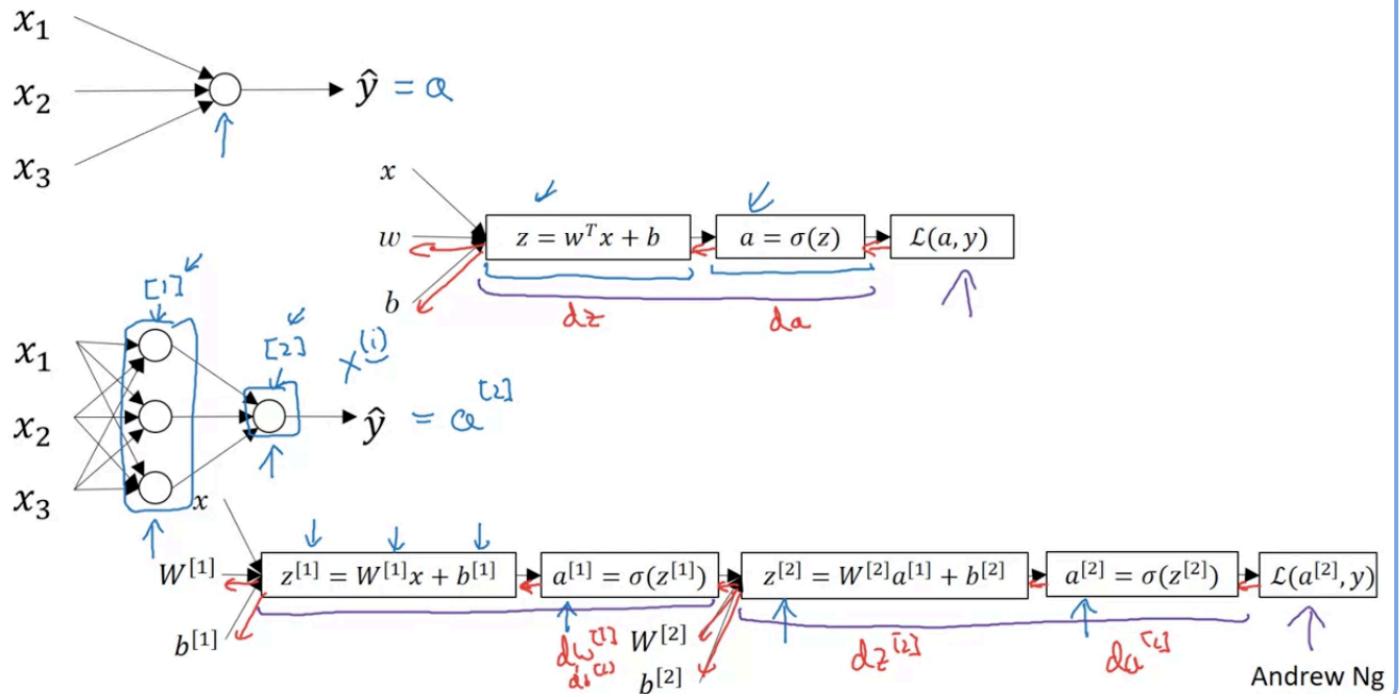
Andrew Ng

Vectorizing Logistic Regression & Gradient Descent...

Broadcasting in Python

## NN Architecture

### What is a Neural Network?

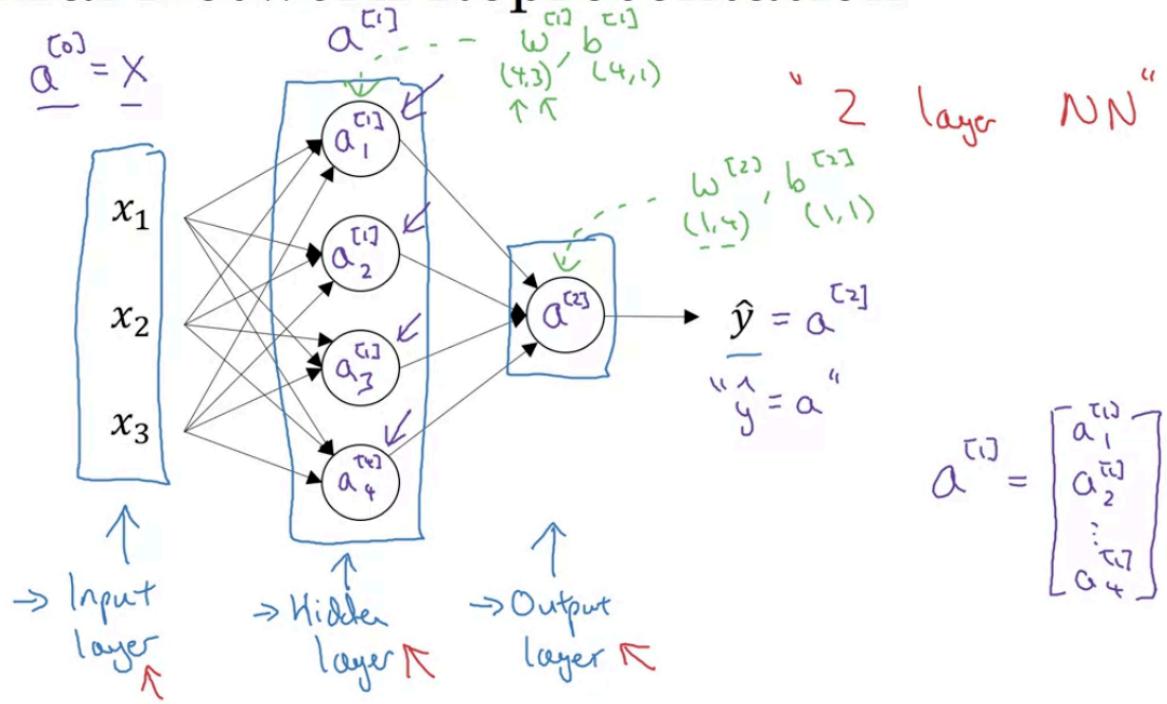


Andrew Ng

Single Hidden Layer example

Input Layer → Hidden Layer → Output Layer

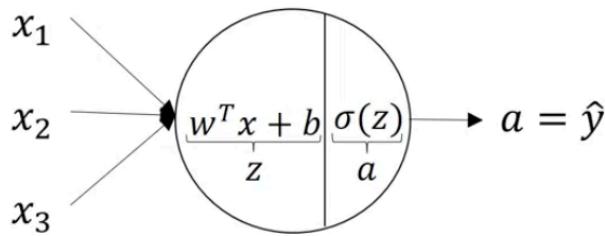
# Neural Network Representation



Andrew Ng

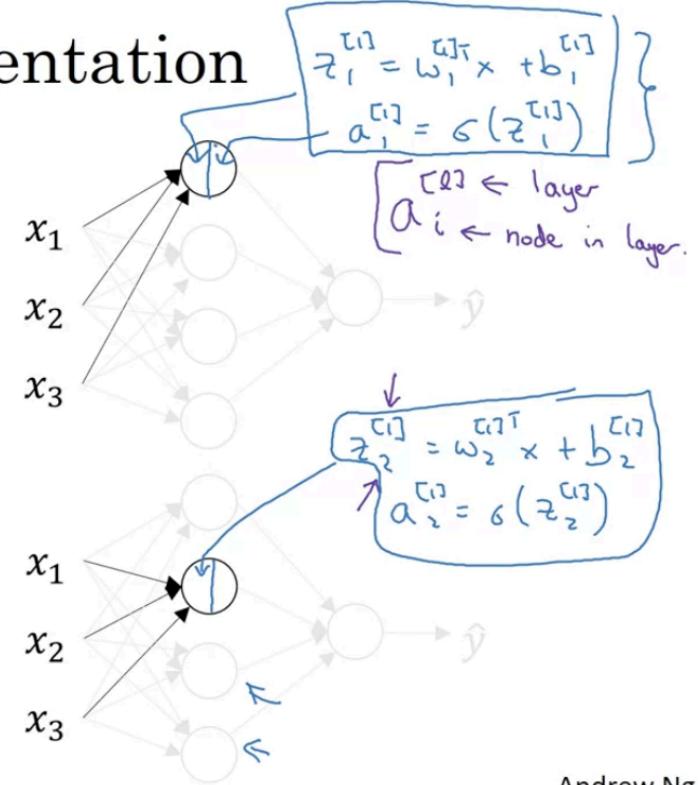
input layer is not counted when referencing the # of layers in NN

# Neural Network Representation



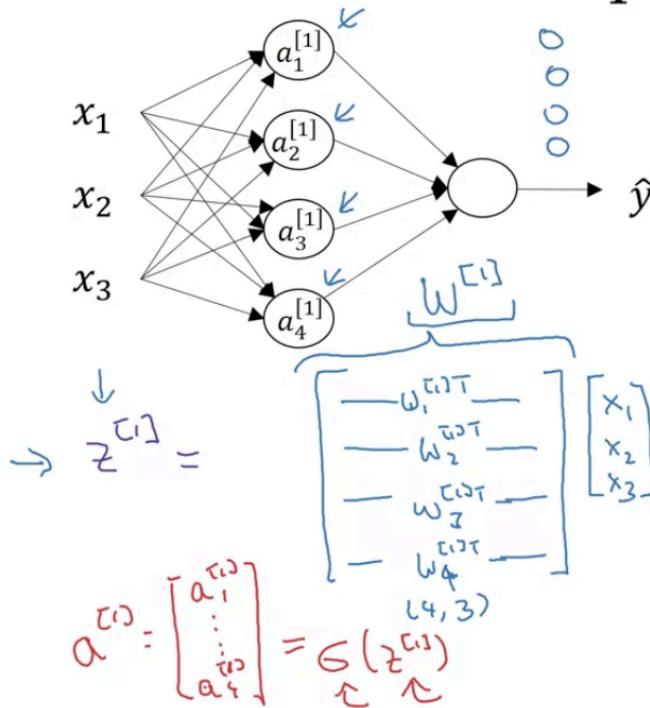
$$z = w^T x + b$$

$$a = \sigma(z)$$



Andrew Ng

# Neural Network Representation



Equation block showing the forward pass of a neural network:

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

$$\rightarrow w_1^{[1]T} x + b_1^{[1]} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

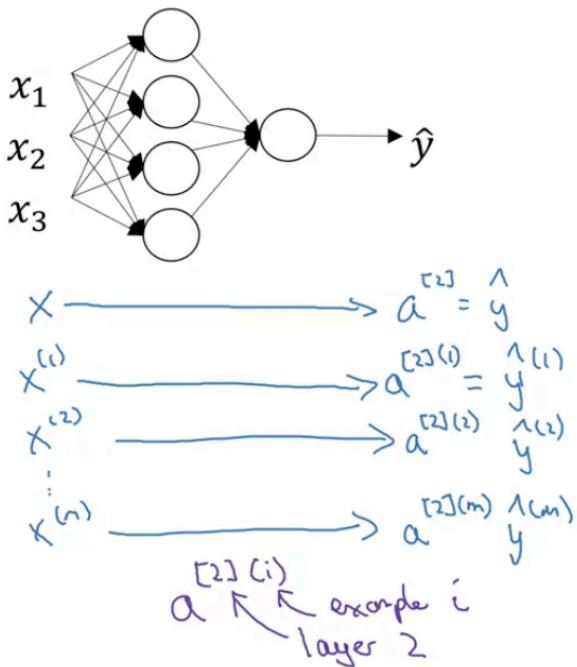
$$\rightarrow w_2^{[1]T} x + b_2^{[1]} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$$\rightarrow w_3^{[1]T} x + b_3^{[1]} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$$\rightarrow w_4^{[1]T} x + b_4^{[1]} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

Andrew Ng

# Vectorizing across multiple examples



$$\left\{ \begin{array}{l} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{array} \right\} \quad \leftarrow$$

for  $i = 1$  to  $n$ ,

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

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

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

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

Andrew Ng

## Vectorizing across multiple examples

for i = 1 to m:

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

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

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

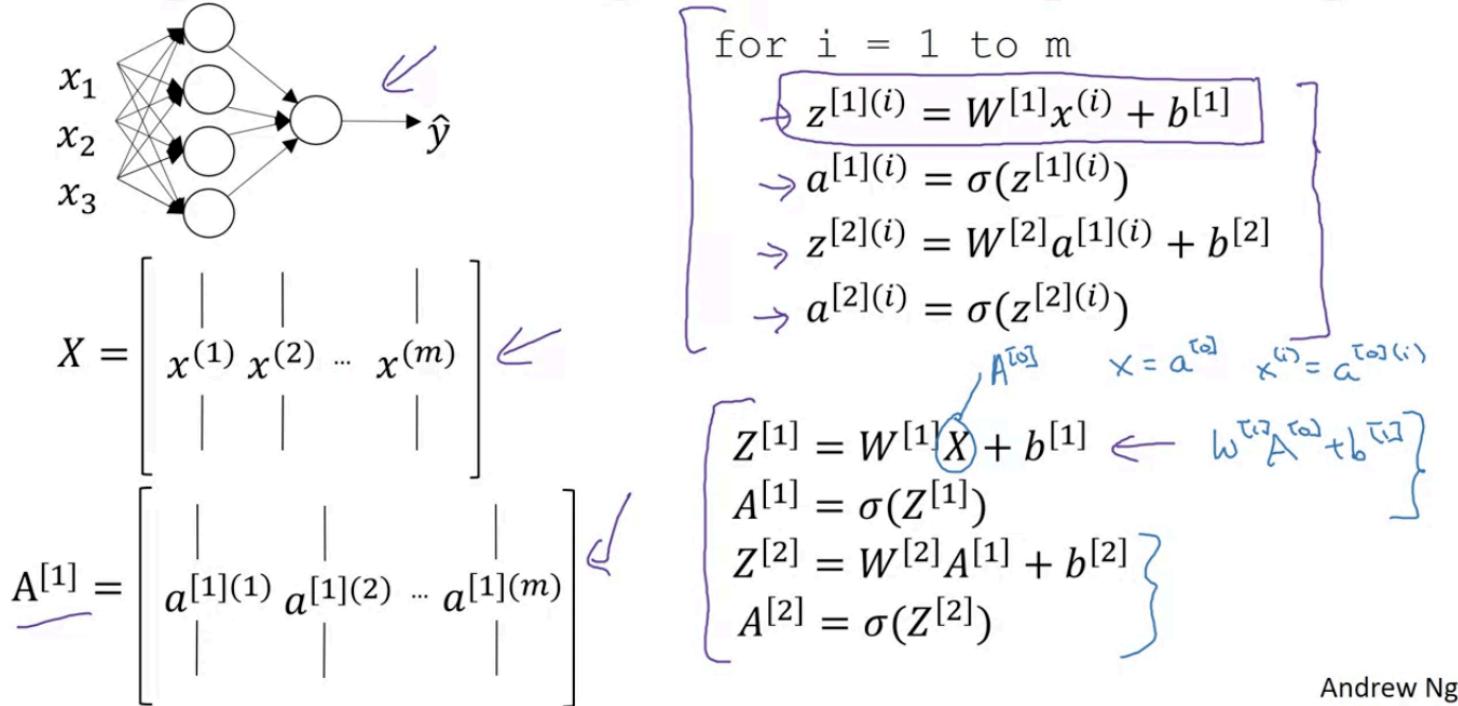
$$a^{[2](i)} = \sigma(z^{[2](i)})$$

$$X = \begin{bmatrix} & & & \\ \text{---} & \text{---} & \text{---} & \text{---} \\ X^{(1)} & X^{(2)} & \dots & X^{(n)} \\ | & | & & | \\ (n_x, n_h) & & & \end{bmatrix}$$

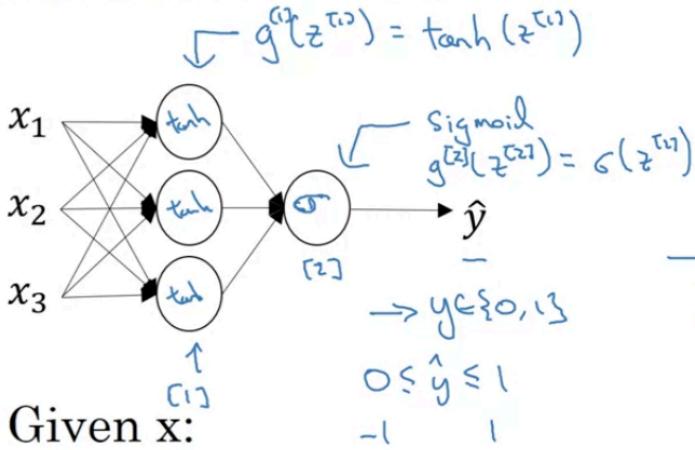
truly complex  
hidden units.

Andrew Ng

# Recap of vectorizing across multiple examples

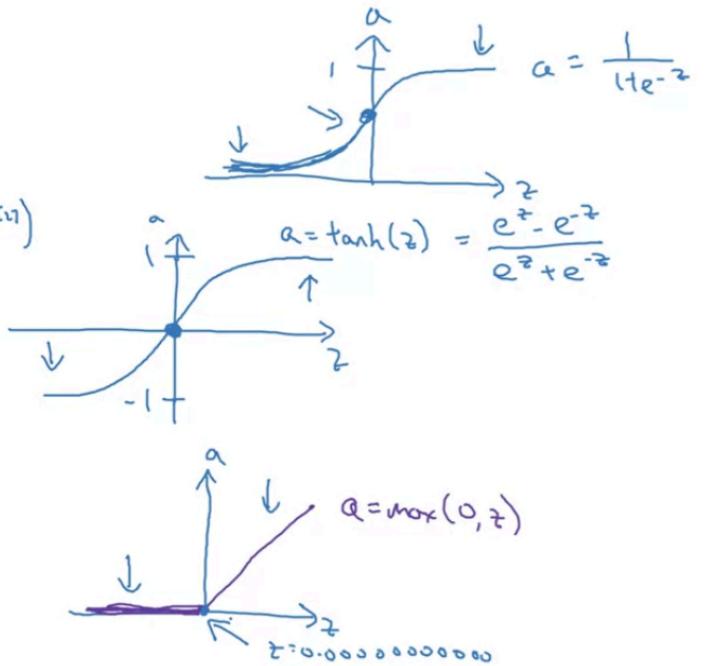


## Activation functions



Given  $x$ :

$$\begin{aligned} z^{[1]} &= W^{[1]}x + b^{[1]} \\ \rightarrow a^{[1]} &= \sigma(z^{[1]}) \quad g^{[1]}(z^{[1]}) \\ z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\ \rightarrow a^{[2]} &= \sigma(z^{[2]}) \quad g^{[2]}(z^{[2]}) \end{aligned}$$



Andrew Ng

Activation Function almost always works better than the sigmoid function

since mean of -1, 1 mostly always becomes nearer to 0

tanh is superior of activation, not practical to use activation nearly for all cases

0-1 value: sigmoid (never use)

ReLU: default for most cases

downside: derivative is 0, when value is negative

to neglect this downside: Leaky ReLU

## Why do you need Non-Linear Activation Functions?

linear hidden layer is useless,  
it is just a logistic regression

Non-linear activation functions are crucial for neural networks because they enable the network to compute more complex and interesting functions. Without them, the network's output would remain a linear function of the input, regardless of the number of hidden layers.

### 1. Linear Activation Limitation:

If the activation function is linear (or identity, where  $g(z) = z$ ), the neural network simply performs linear operations. This results in the output being just a linear transformation of the input:

$$a^{(1)} = W^{(1)}x + b^{(1)}$$
$$a^{(2)} = W^{(2)}a^{(1)} + b^{(2)} = W^{(\text{prime})}x + b^{(\text{prime})}$$

This makes the hidden layers redundant, as the entire network simplifies to a linear model.

### 2. Deep Networks and Linear Activation:

Even with many hidden layers, if all activations are linear, the overall network still computes a linear function. Therefore, the depth of the network does not increase its expressive power. The text emphasizes that without non-linearity, you might as well not have hidden layers.

### 3. Expressiveness of Non-Linear Activation:

Non-linear activation functions (like **ReLU**, **tanh**, **sigmoid**) allow the network to compute more complex, non-linear functions. This is essential for learning complicated patterns and relationships in data.

### 4. Comparison to Logistic Regression:

A neural network with linear hidden layers and a sigmoid output is no more powerful than logistic regression. This highlights the importance of non-linearity in hidden layers to enhance the network's capacity.

### 5. When Linear Activation is Useful:

- **Output Layer for Regression:** Linear activation is useful in the output layer for regression tasks, where the target output is a continuous value (e.g., predicting housing prices).
- For hidden layers, non-linear functions are preferred to introduce the complexity needed for learning diverse patterns.

## Derivatives of Activation Functions - Back Propagation

*Sigmoid*

*tanh*

*ReLU and Leaky ReLU*

## Gradient Descent for NNs

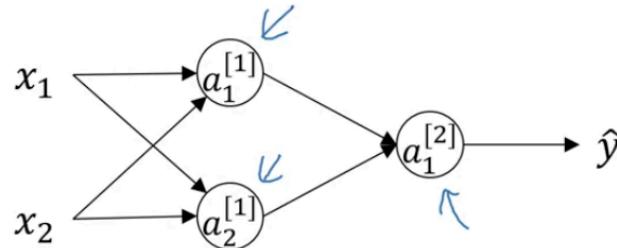
for gradient descent start with random values for parameters

## Random Initialization

when started with zero, different neurons are initialized in the same way and compute the same operations.

How big?

## Random initialization



$$\begin{aligned} \rightarrow w^{[1]} &= \text{np. random. randn}((2, 2)) * \frac{0.01}{100?} & \rightarrow z^{[1]} &= w^{[1]} x + b^{[1]} \\ b^{[1]} &= \text{np. zeros}((2, 1)) & a^{[1]} &= g^{[1]}(z^{[1]}) \\ w^{[2]} &= \text{np. random. randn}((1, 2)) * 0.01 & \downarrow & \\ b^{[2]} &= 0 & \uparrow & \end{aligned}$$

A graph of the sigmoid function  $g(z)$ . The horizontal axis is labeled  $z$  and the vertical axis is labeled  $g(z)$ . The curve is S-shaped, passing through the point (0, 0.5).

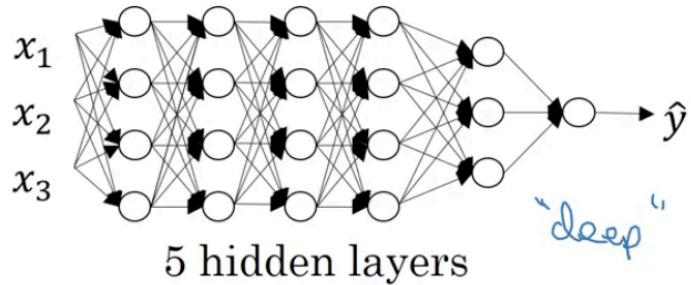
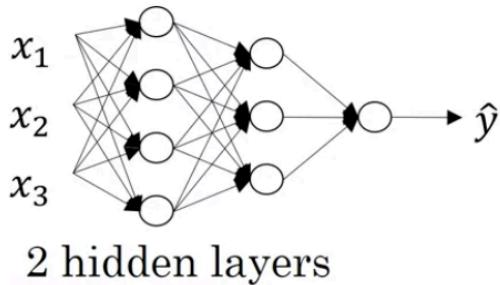
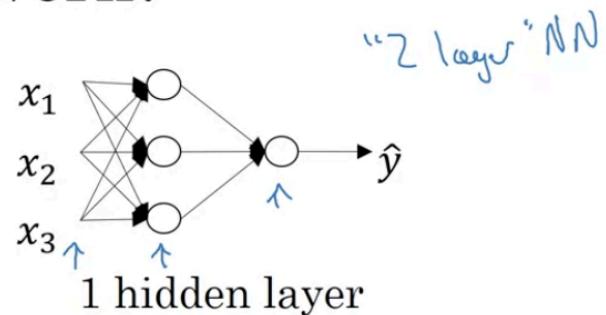
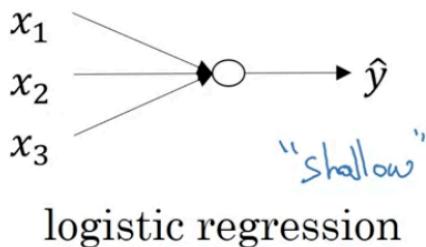
Andrew Ng

when  $w$  is too large, initialized at plateau points

## Deep L-layer NN

Deep NN ?

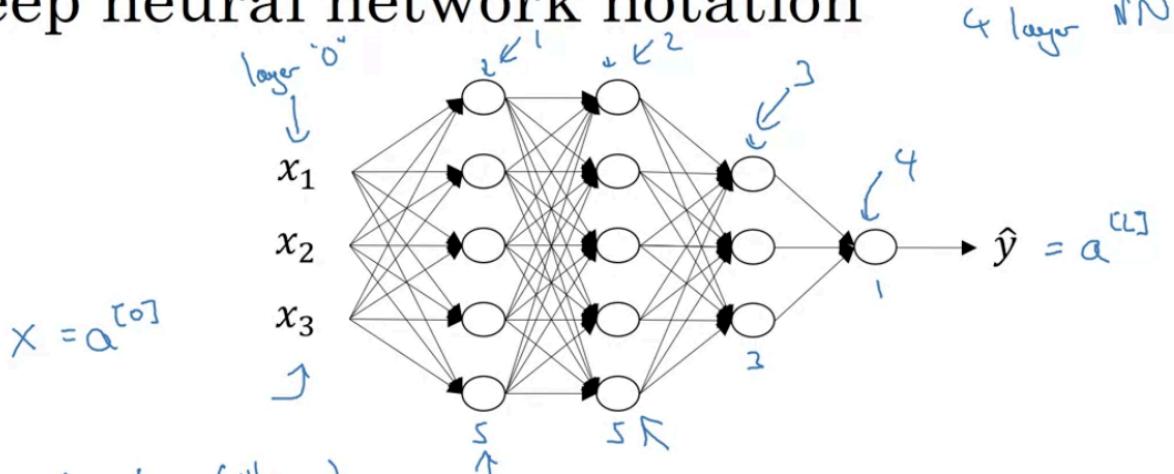
# What is a deep neural network?



Andrew Ng

# of hidden layers is another 'hyperparameter' for some cases

## Deep neural network notation



$$l = 4 \quad (\# \text{layers})$$

$$n^{[l]} = \# \text{units in layer } l$$

$$a^{[l]} = \text{activations in layer } l$$

$$a^{[l]} = g^{[l]}(z^{[l]}), \quad w_{j,l}^{[l]} = \text{weights for } z_j^{[l]}$$

$$n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3, n^{[4]} = n^{[L]} = 1$$

$$n^{[0]} = n_x = 3$$

Andrew Ng

- Layer "0" (input layer)
- $x_1, x_2, x_3$  (input features)
- $X = a^{[0]}$
- 4 layer NN (neural network)
- $l = 4$  (# layers)
- $n^{[l]}$  = # units in layer  $l$
- $a^{[l]}$  = activations in layer  $l$
- $a^{[l]} = g(z^{[l]})$

- $W^{[l]}$ ,  $b^{[l]}$  = weights for  $z^{[l]}$
- $n^{[1]} = 5$ ,  $n^{[2]} = 5$ ,  $n^{[3]} = 3$ ,  $n^{[4]} = 1$
- $n^{[0]} = n_x = 3$

### Neural Network Visualization (Text Representation):

#### • Layer 0 (Input Layer):

- 3 input features  $x_1, x_2, x_3$
- $n^{[0]} = 3$  (number of neurons)

#### • Layer 1 (Hidden Layer 1):

- 5 neurons
- $n^{[1]} = 5$

#### • Layer 2 (Hidden Layer 2):

- 5 neurons
- $n^{[2]} = 5$

#### • Layer 3 (Hidden Layer 3):

- 3 neurons
- $n^{[3]} = 3$

#### • Layer 4 (Output Layer):

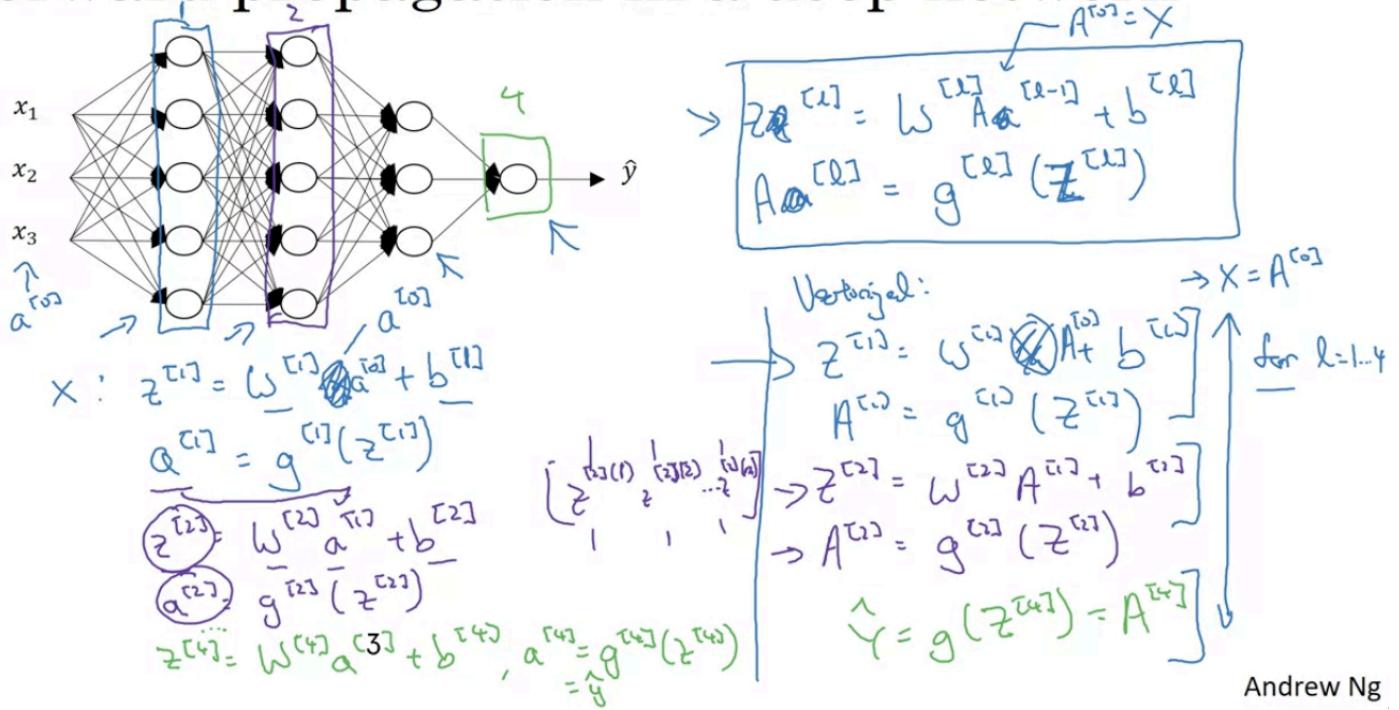
- 1 neuron (output)
- $n^{[4]} = 1$

#### • Total Layers ( $l$ ):

- $l = 4$

## Forward Propagation in a DNN

### Forward propagation in a deep network

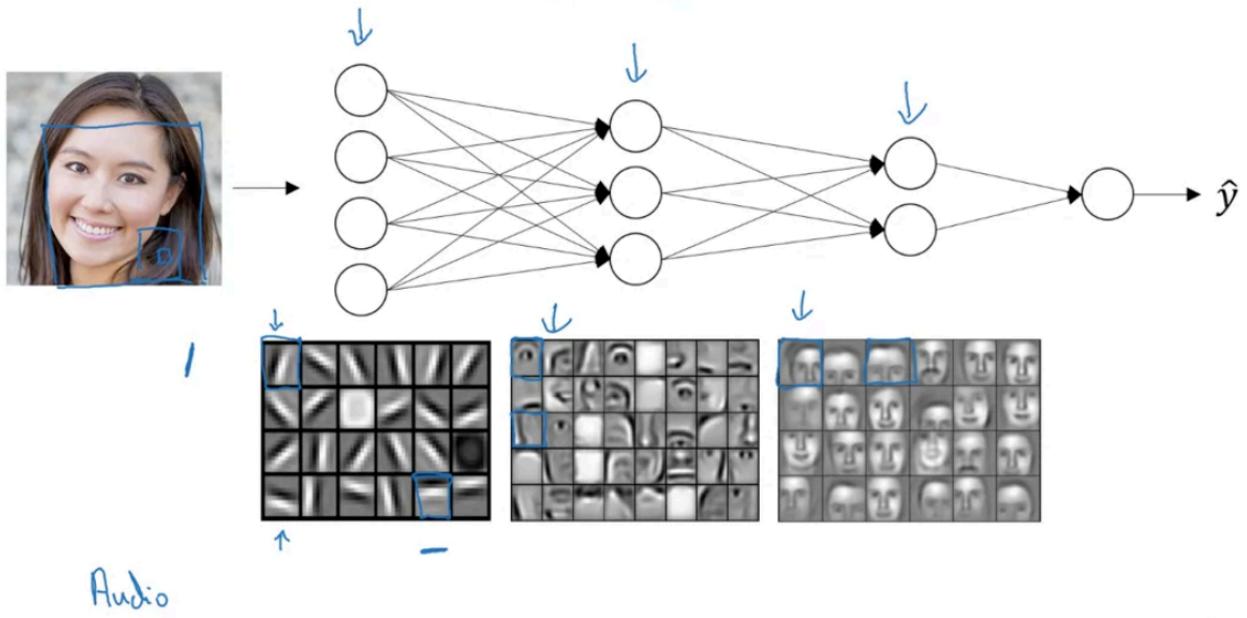


Get your Matrix Dimensions right

Why deep networks?

that famous representation

# Intuition about deep representation

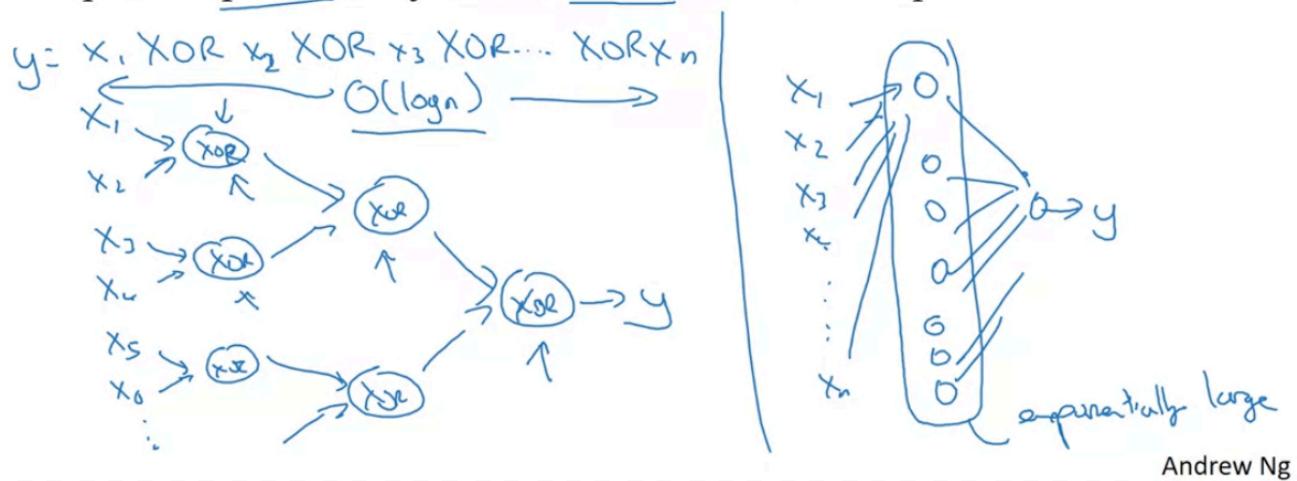


ex: for audio - low level - audio waveform (different sounds) - words - sentences/phrases

## Circuit Theory

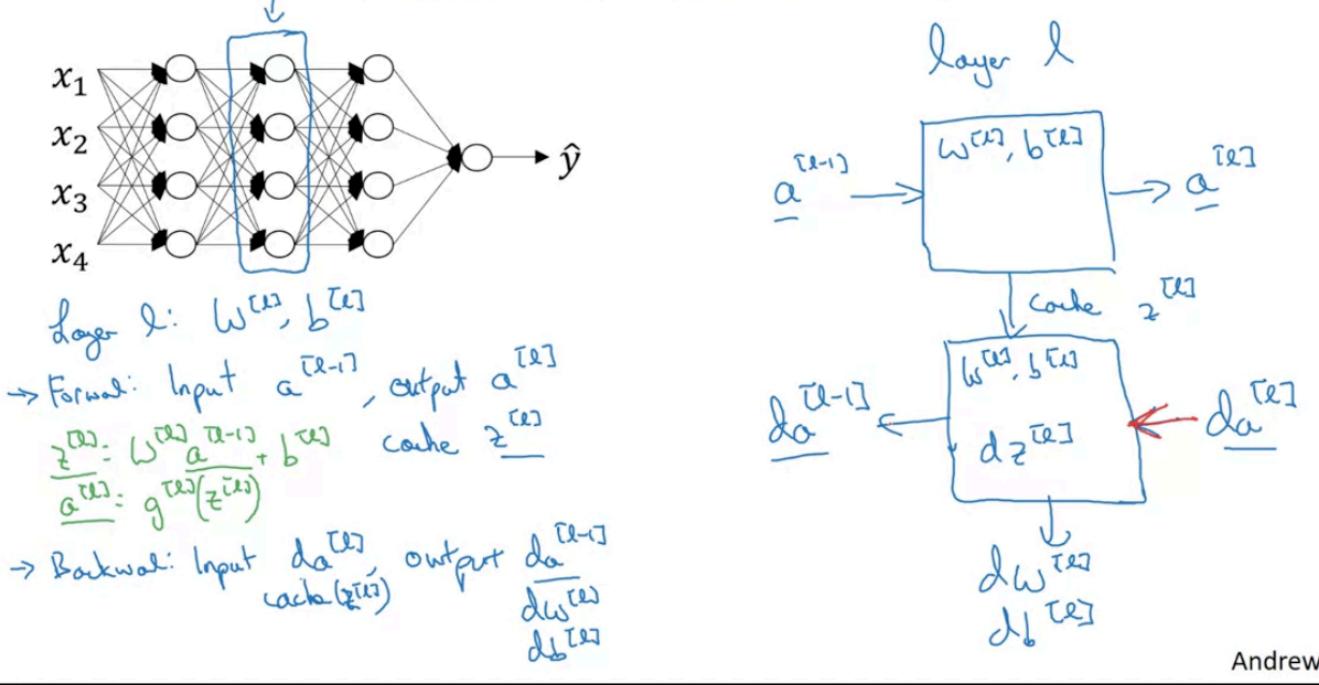
## Circuit theory and deep learning

Informally: There are functions you can compute with a “small” L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

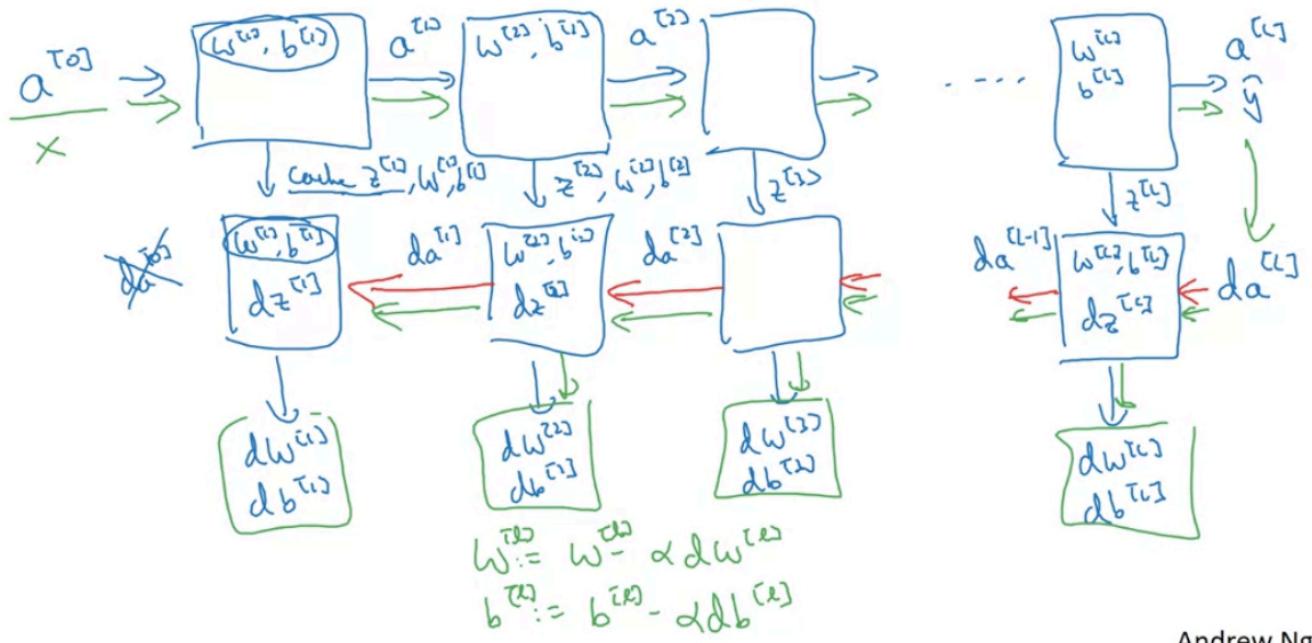


## Building blocks of Deep NNs:

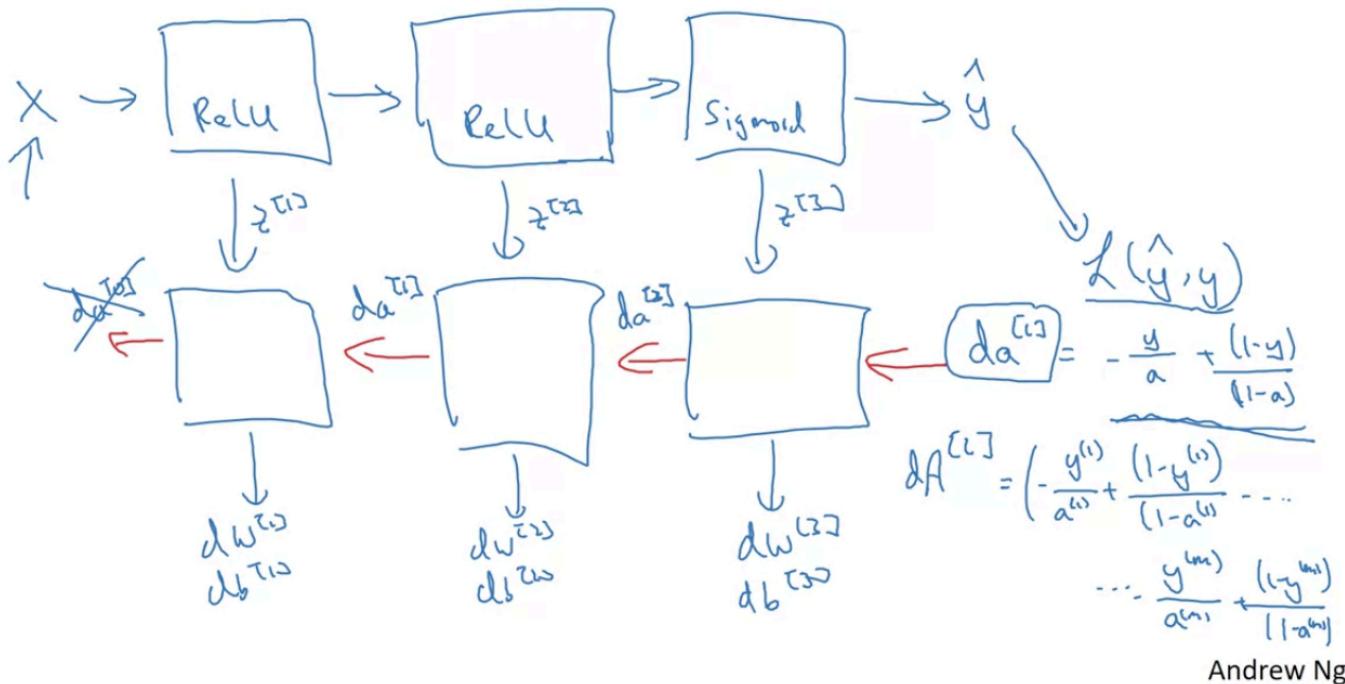
# Forward and backward functions



# Forward and backward functions



# Summary



## Parameters & Hyperparameters ?

What are hyperparameters?

Parameters:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]} \dots$

Hyperparameters: learning rate  $\alpha$   
#iterations  
#hidden layers  $L$   
#hidden units  $n^{[1]}, n^{[2]}, \dots$   
choice of activation function

Later: Momentum, mini-batch size, regularizations, ...

Andrew Ng

## Deep Learning and Brain ?

much more complex, can not even understand a single neuron yet  
it just *resembles*

## Summary of Neural Network Training and Inference (Running the Model):

### 1. Training a Neural Network (Step-by-Step):

Training involves **forward propagation**, **loss computation**, **backward propagation**, and **parameter update**. This process is repeated for multiple iterations (epochs) to minimize the cost function.

## Step 1: Forward Propagation (Compute Predictions)

Purpose: Calculate predictions by passing input data through the network layer by layer.

### Steps (per layer):

#### 1. Linear Step:

- Linear output before activation
- Weights for layer
- Activations from the previous layer (or input for the first layer)
- Bias vector for layer

#### 2. Activation Step:

- is the activation function:
- **ReLU** for hidden layers:
- **Sigmoid** for the output layer (for binary classification):

## Step 2: Compute Cost (Loss Function)

Purpose: Measure how far predictions are from the true labels.

### Cost Function (Binary Cross-Entropy):

- Number of examples
- True labels
- Predicted output from the final layer

## Step 3: Backward Propagation (Compute Gradients)

Purpose: Calculate gradients to adjust weights and biases to minimize the loss.

### Steps (per layer, from last to first):

#### 1. Output Layer Gradient (Final Layer):

- Gradient of the loss with respect to .

#### 2. Hidden Layer Gradients (for each layer from to 1):

- Derivative of the activation function (ReLU or sigmoid).

#### 3. Gradients for Weights and Biases:

## Step 4: Update Parameters (Gradient Descent Step)

Purpose: Adjust weights and biases using gradients to minimize the cost.

- Learning rate (controls the size of updates).

## 2. Running the Neural Network (After Training - Inference):

Now that the model is trained, inference involves **forward propagation only**.

## **Steps for Inference (Forward Propagation):**

1. **Input Data** – Pass input through the network.
2. **Linear and Activation Steps (Layer by Layer)**:
  - is **ReLU** for hidden layers and **sigmoid** for the final layer.
3. **Final Output**:
  - – Prediction (probability between 0 and 1 for binary classification).
4. **Thresholding for Binary Classification**:

## **Order of Operations (Training Phase):**

1. **Forward Propagation (Layer by Layer)**:
  - Linear Step:
  - Activation:
2. **Compute Cost**
3. **Backward Propagation (Layer by Layer, Reverse Order)**:
  - Compute for output layer
  - Calculate gradients
4. **Parameter Update**:
5. **Repeat for Multiple Iterations**

## **Order of Operations (Inference Phase - Running the Model):**

1. **Input Data**:
2. **Forward Propagation (Layer by Layer)**:
3. **Final Prediction**:
  - (Output Layer)
4. **Thresholding for Binary Classification**:

## **Key Concepts Simplified:**

- **Forward Propagation**: Pass input through the network to compute predictions.
- **Cost Function**: Measures how good/bad the predictions are.
- **Backward Propagation**: Calculate gradients to improve weights/biases.
- **Parameter Update**: Adjust weights/biases to reduce the cost.
- **Repeat**: Continue this process until the model converges (low cost).