# Artificial Intelligence and Machine Learning
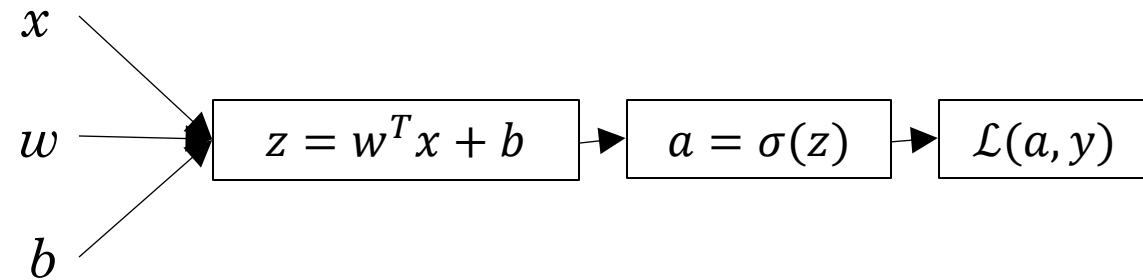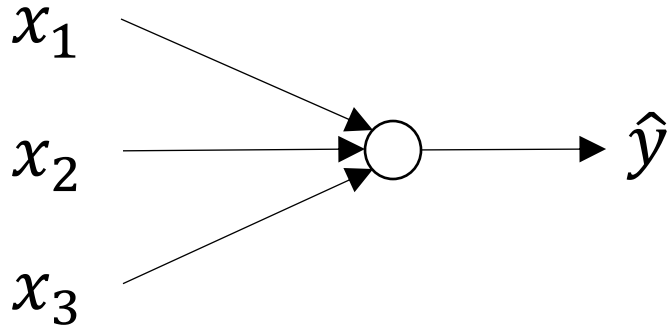
# Neural Networks

# Lecture Outline

- Logistic Regression Review
- Neural Networks
  - Forward pass
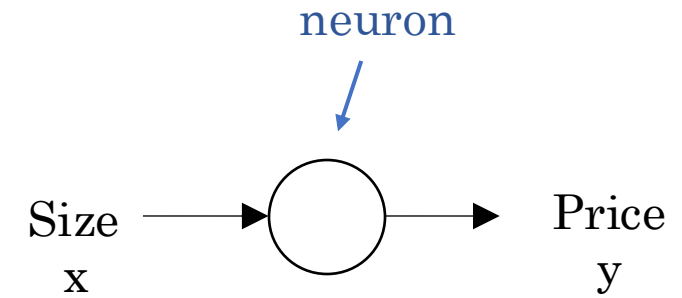  - Backward pass

# Review: Logistic Regression
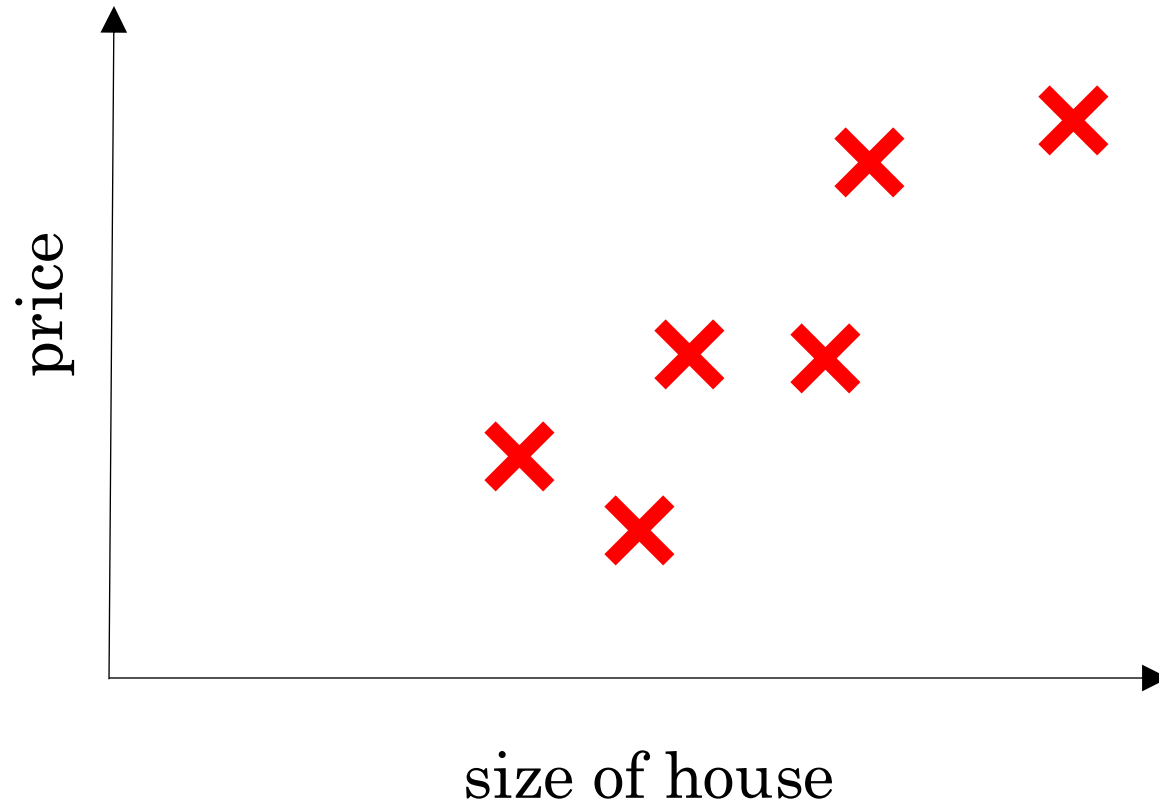
$x_1$

$x_2$     ○  →   $\hat{y}$

$x_3$

$x$

$w$   ◆   $\boxed{z = w^T x + b}$ → $\boxed{a = \sigma(z)}$ → $\boxed{\mathcal{L}(a, y)}$

$b$

# Introduction to Deep Learning

# What is a Neural Network?

deeplearning.ai

# Housing Price Prediction

price

size of house

neuron

Size
x → ◯ → Price
y

# Housing Price Prediction

size

#bedrooms ───▶ ◯ family size

zip code ───▶ ◯ walkability

wealth ───▶ ◯ school quality

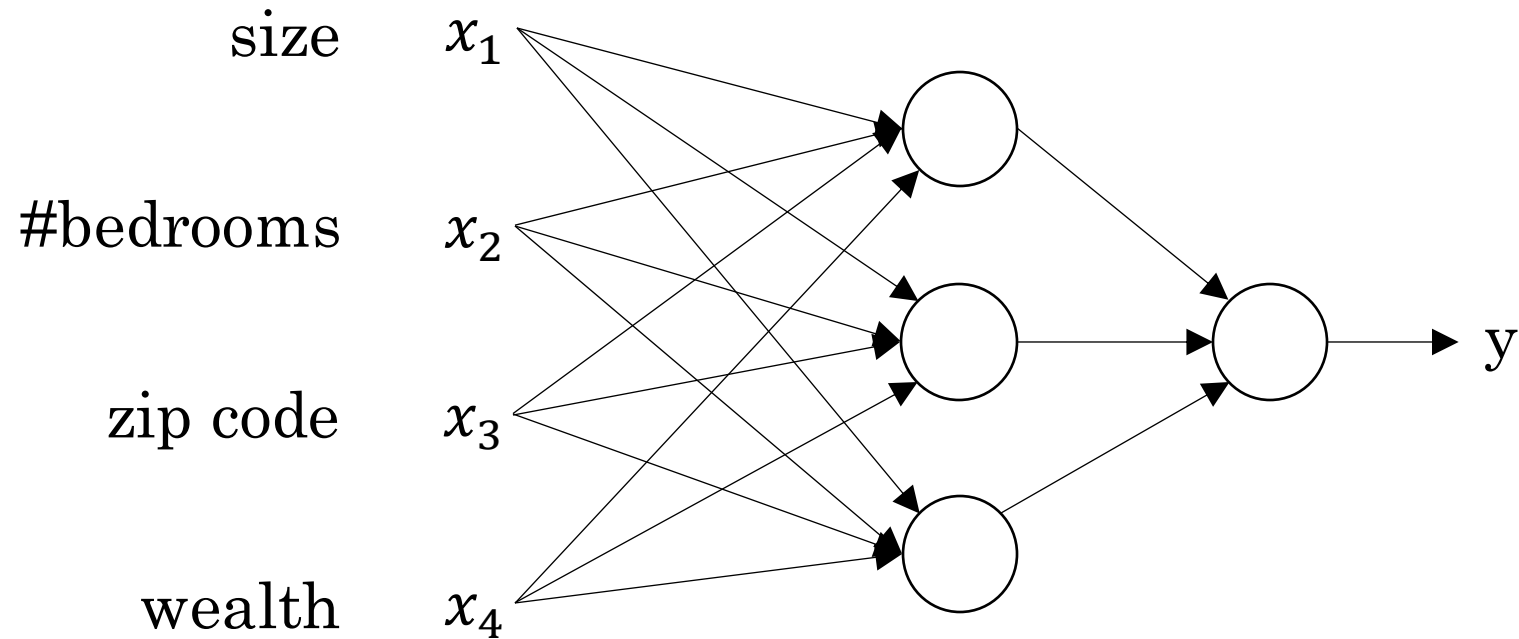# Housing Price Prediction

# Housing Price Prediction

# Introduction to Deep Learning

## Supervised Learning with Neural Networks

deeplearning.ai

# Supervised Learning

| Input(x) | Output (y) | Application |
|---|---|---|
| Home features | Price | Real Estate |
| Ad, user info | Click on ad? (0/1) | Online Advertising |
| Image | Object (1,…,1000) | Photo tagging |
| Audio | Text transcript | Speech recognition |
| English | Chinese | Machine translation |
| Image, Radar info | Position of other cars | Autonomous driving |

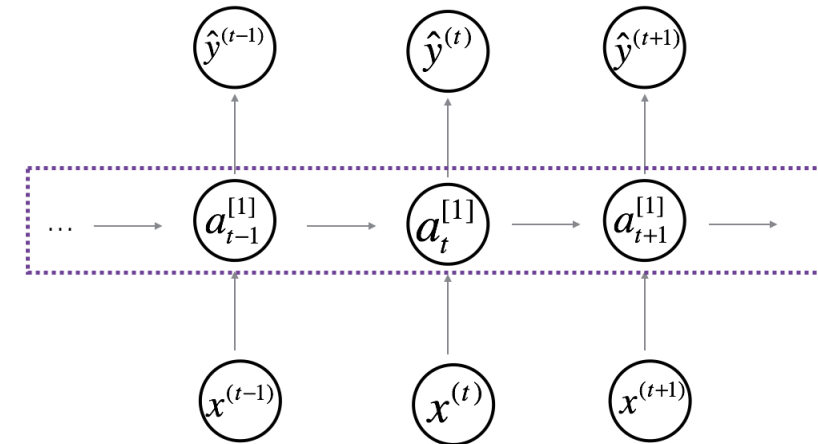# Neural Network examples
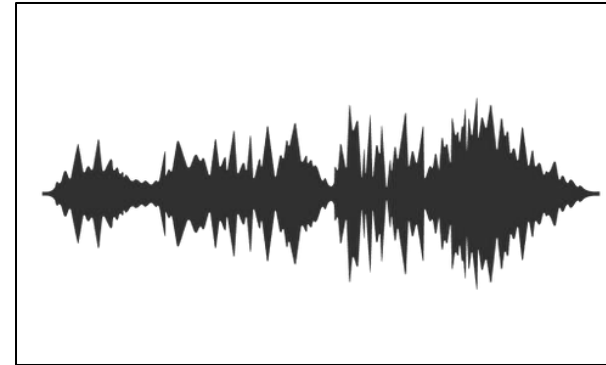


Standard NN         Convolutional NN         Recurrent NN

# Supervised Learning

## Structured data

| Size | #bedrooms | ... | Price (1000$s) |
|------|-----------|-----|----------------|
| 2104 | 3 | | 400 |
| 1600 | 3 | | 330 |
| 2400 | 3 | | 369 |
| ⋮ | ⋮ | | ⋮ |
| 3000 | 4 | | 540 |

| User Age | Ad ID | ... | Click |
|----------|-------|-----|-------|
| 41 | 93242 | | 1 |
| 80 | 93287 | | 0 |
| 18 | 87312 | | 1 |
| ⋮ | ⋮ | | ⋮ |
| 27 | 71244 | | 1 |

## Unstructured data



Audio



Image

Four score and seven years ago

Text

One hidden layer
Neural Network

Neural Networks
Overview

deeplearning.ai

# What is a Neural Network?



$x_1$
$x_2$
$x_3$
$\hat{y}$

$x$
$w$
$b$

$$z = w^T x + b \rightarrow a = \sigma(z) \rightarrow \mathcal{L}(a, y)$$

$x_1$
$x_2$
$x_3$
$\hat{y}$

$x$
$W^{[1]}$
$b^{[1]}$

$$z^{[1]} = W^{[1]}x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow \mathcal{L}(a^{[2]}, y)$$

$W^{[2]}$
$b^{[2]}$

One hidden layer
Neural Network

Neural Network
Representation

deeplearning.ai

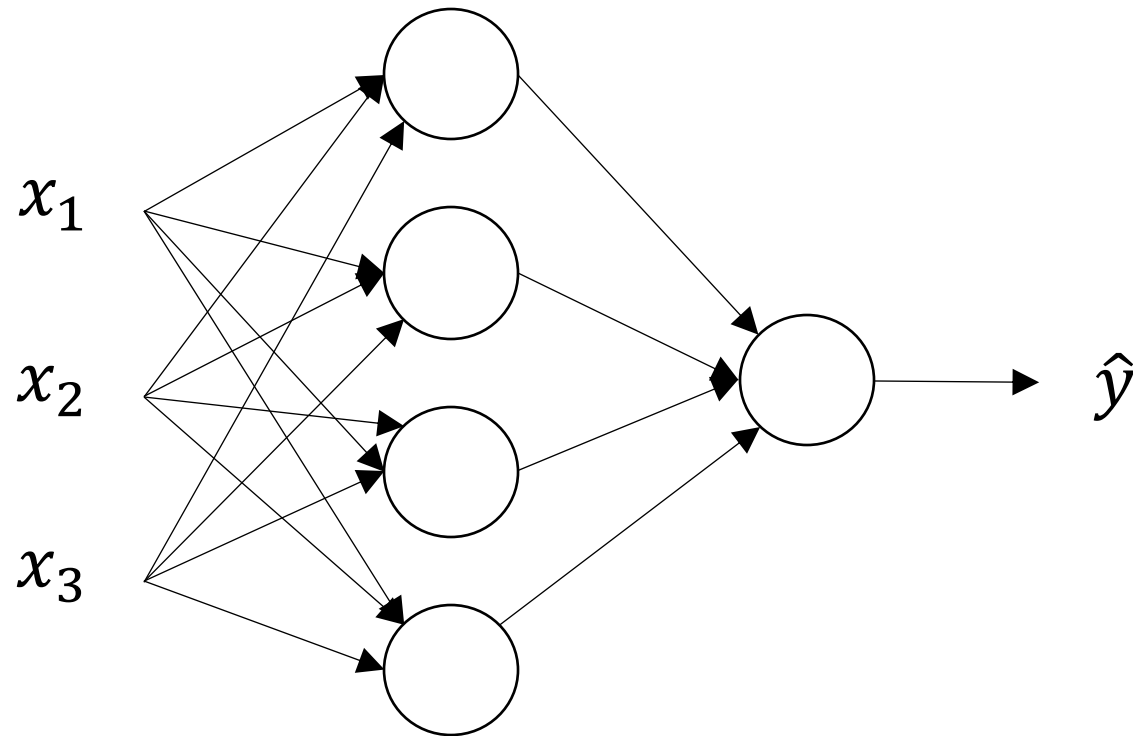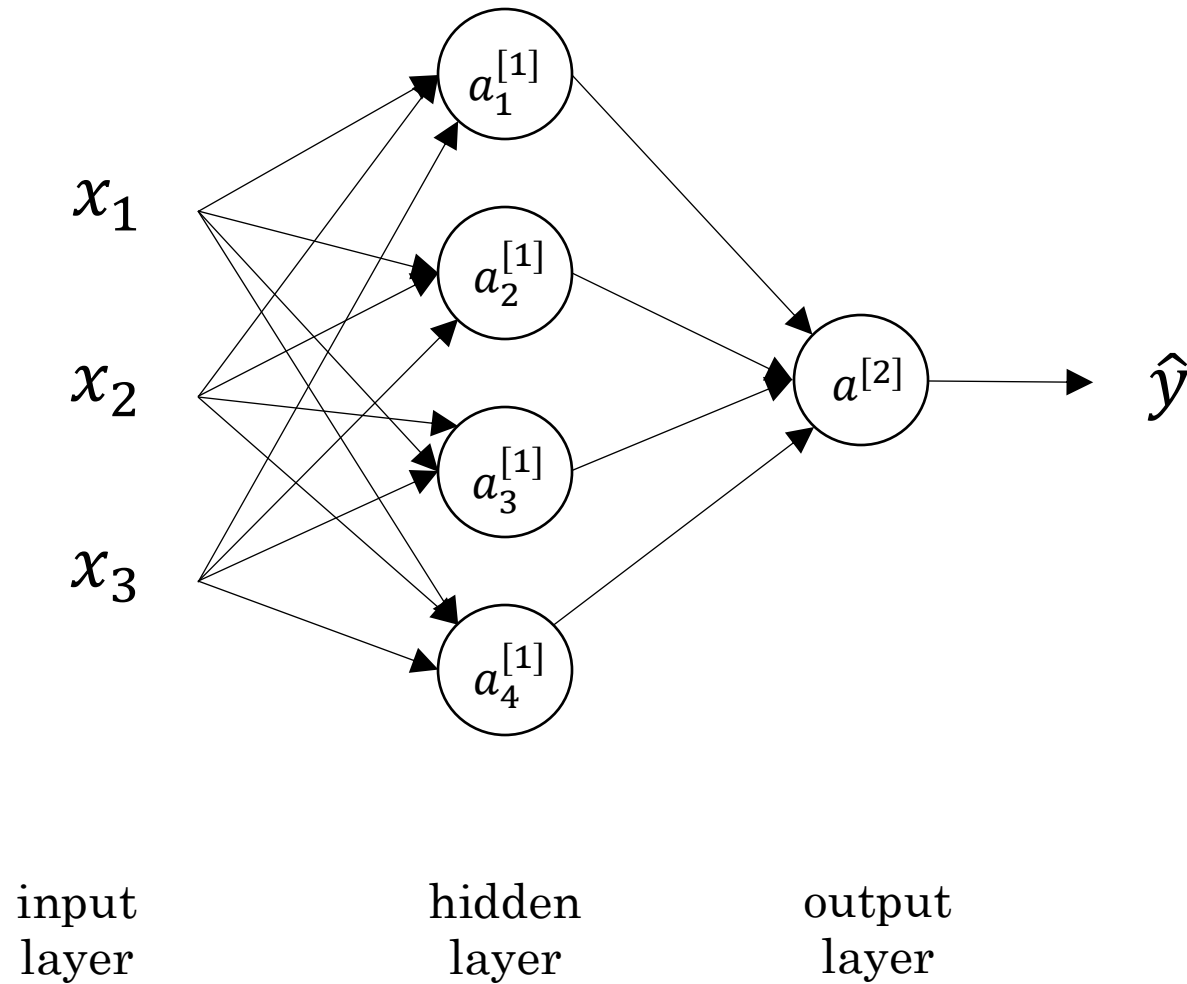# Neural Network Representation

# Neural Network Representation

deeplearning.ai

# One hidden layer Neural Network

## Computing a Neural Network's Output

# Neural Network Representation



$x_1$
$x_2$
$x_3$

$\underbrace{w^T x + b}_{z} \Big| \underbrace{\sigma(z)}_{a}$
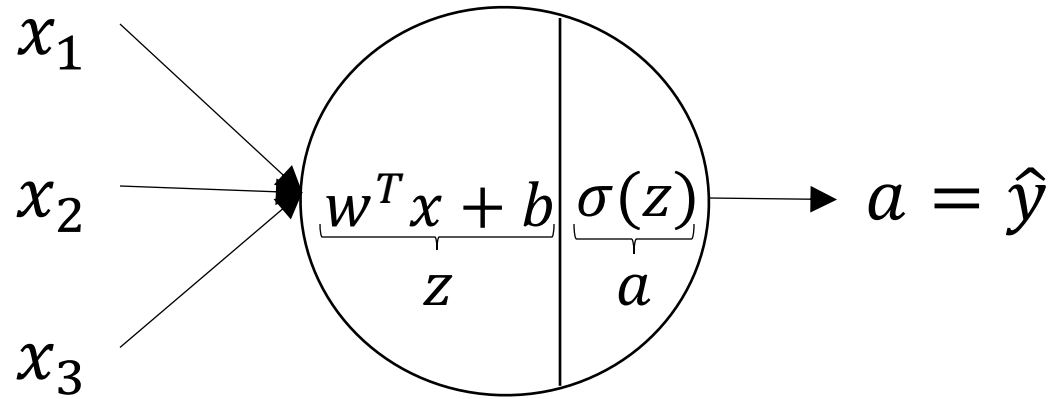
$a = \hat{y}$

$z = w^T x + b$

$a = \sigma(z)$

$x_1$
$x_2$
$x_3$

$\hat{y}$

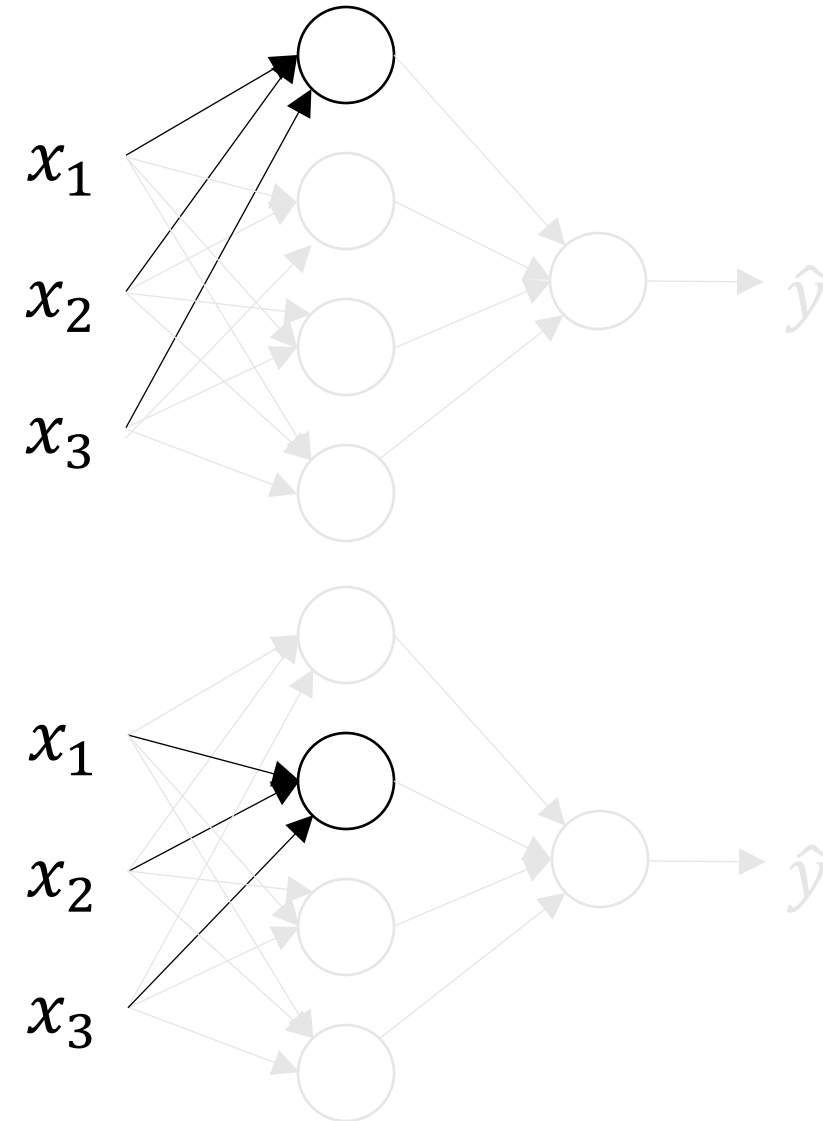# Neural Network Representation



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

$$a = \sigma(z)$$

# Neural Network Representation

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \ a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \ a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \ a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \ a_4^{[1]} = \sigma(z_4^{[1]})$$

# Neural Network Representation learning



Given input x:

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

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

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

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

# One hidden layer Neural Network

deeplearning.ai

# Vectorizing across multiple examples

# Vectorizing across multiple examples



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

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

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

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

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

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

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

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

# Vectorizing across multiple examples



$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} | & | & & | \\ a^{[1](1)} & a^{[1](2)} & \cdots & a^{[1](m)} \\ | & | & & | \end{bmatrix}$$

```
for i = 1 to m
```

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

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$
$$A^{[1]} = \sigma(Z^{[1]})$$
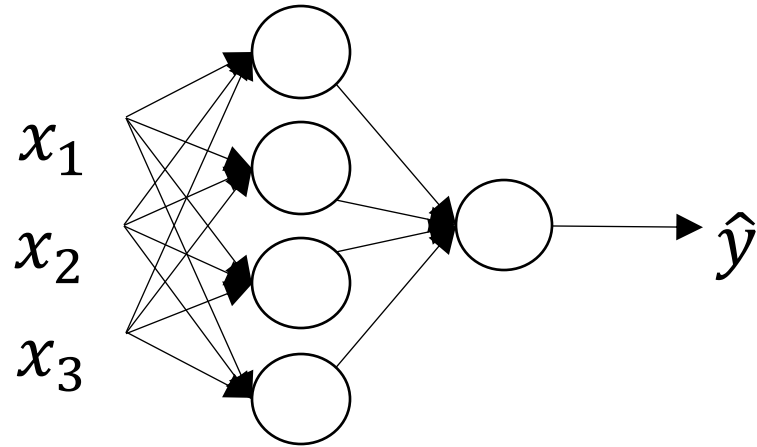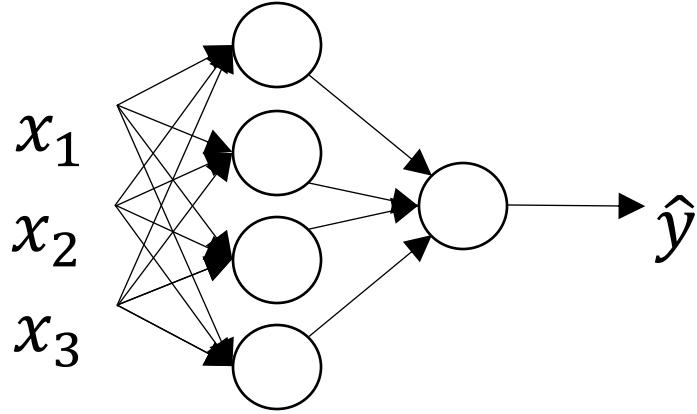$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$
$$A^{[2]} = \sigma(Z^{[2]})$$

One hidden layer
Neural Network

Explanation
for vectorized
implementation

deeplearning.ai

# Justification for vectorized implementation

$$z^{[1](1)} = w^{[1]}x^{(1)} + b^{[1]}, \quad z^{[1](2)} = w^{[1]}x^{(2)} + b^{[1]}, \quad z^{[1](3)} = w^{[1]}x^{(3)} + b^{[1]}$$

$$w^{[1]} = \begin{bmatrix} \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \end{bmatrix} \qquad w^{[1]}x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \qquad w^{[1]}x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \qquad w^{[1]}x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

# Justification for vectorized implementation

$$z^{[1](1)} = W^{[1]} x^{(1)} + b^{[1]} \quad, \quad z^{[1](2)} = W^{[1]} x^{(2)} + b^{[1]} \quad, \quad z^{[1](3)} = W^{[1]} x^{(3)} + b^{[1]}$$

$$W^{[1]} = \begin{bmatrix} \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} \end{bmatrix} \qquad W^{[1]} x^{(1)} = \begin{bmatrix} \bullet \\ \bullet \\ \vdots \\ \bullet \end{bmatrix} \qquad W^{[1]} x^{(2)} = \begin{bmatrix} \bullet \\ \bullet \\ \vdots \\ \bullet \end{bmatrix} \qquad W^{[1]} x^{(3)} = \begin{bmatrix} \circ \\ \circ \\ \vdots \\ \circ \end{bmatrix}$$

$$z^{[1]} = W^{[1]} X + b^{[1]} \qquad W^{[1]} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \cdots \\ | & | & | \end{bmatrix} = \begin{bmatrix} \bullet & \bullet & \circ \\ \bullet & \bullet & \circ \\ \bullet & \bullet & \circ \\ \bullet & \bullet & \circ \end{bmatrix} = \begin{bmatrix} | & | & | \\ z^{[1](1)} & z^{[1](2)} & z^{[1](2)} \cdots \\ | & | & | \end{bmatrix} = z^{[1]}$$

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

$+ b^{[1]} \qquad + b^{[1]} \qquad + b^{[1]}$

Andrew Ng

# Recap of vectorizing across multiple examples

$x_1$
$x_2$
$x_3$

$\hat{y}$

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} | & | & & | \\ a^{[1](1)} & a^{[1](2)} & \cdots & a^{[1](m)} \\ | & | & & | \end{bmatrix}$$

```
for i = 1 to m
```

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

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

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

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

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

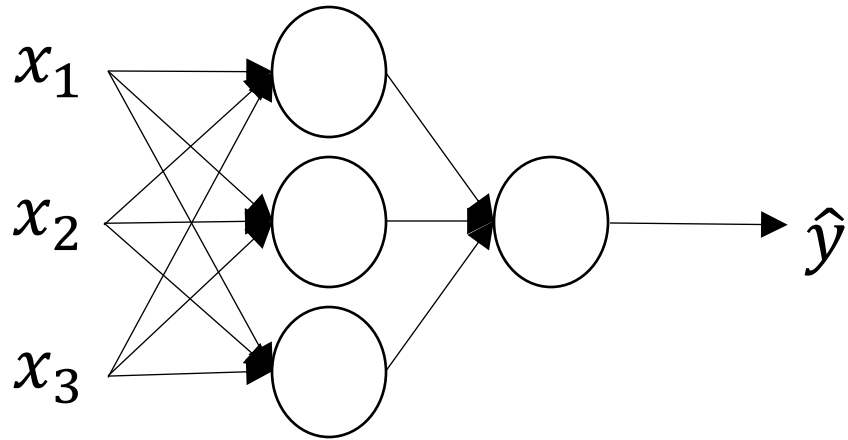$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]})$$

One hidden layer
Neural Network

Activation functions

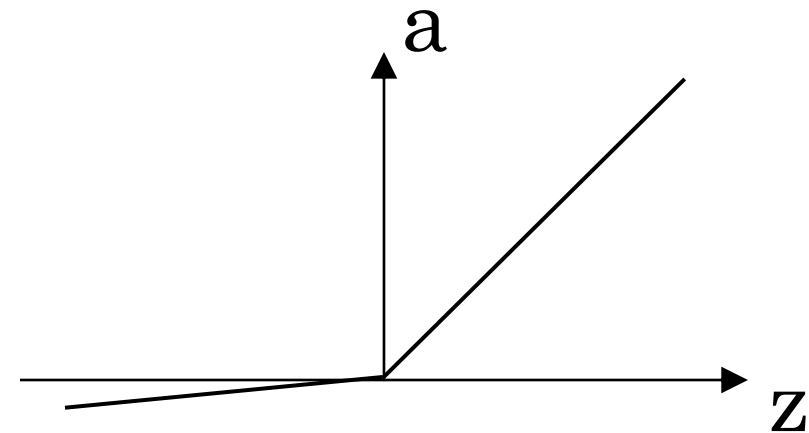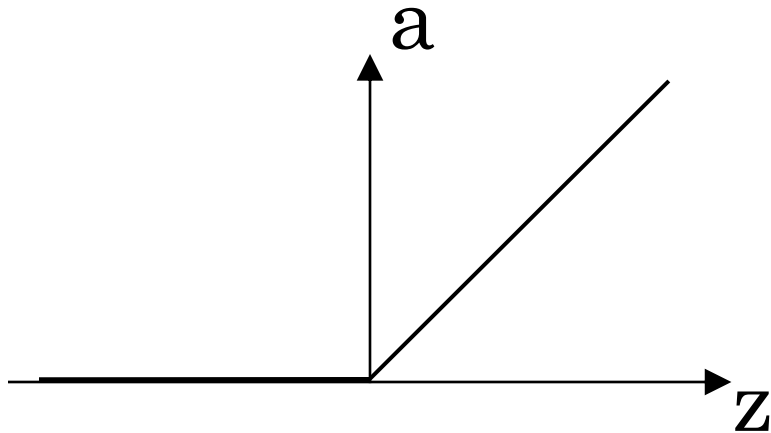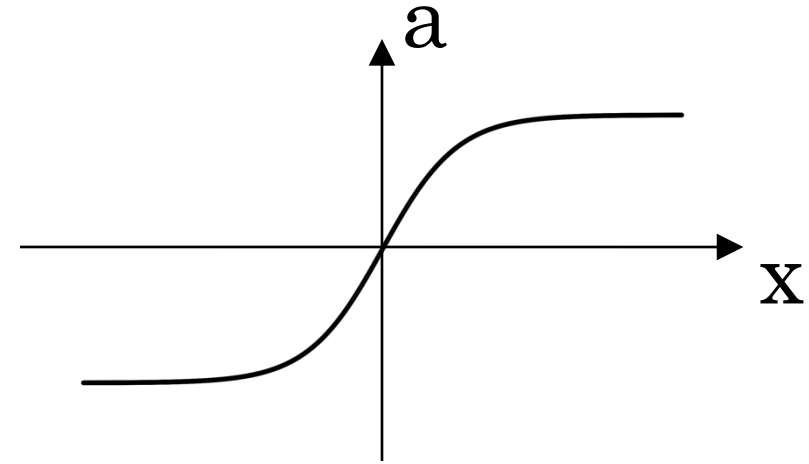deeplearning.ai

# Activation functions



Given x:

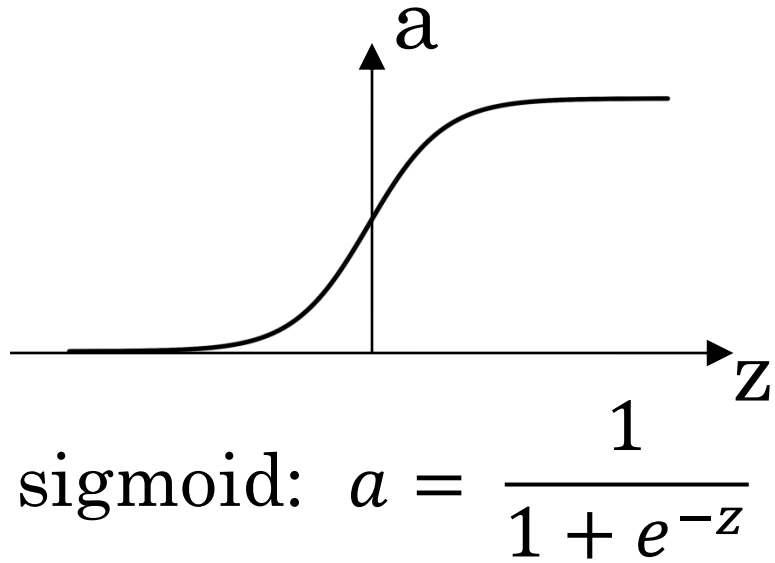$$z^{[1]} = W^{[1]}x + b^{[1]}$$

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

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

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

# Pros and cons of activation functions

sigmoid: $a = \dfrac{1}{1 + e^{-z}}$

# Pros and cons of activation functions



sigmoid: $a = \dfrac{1}{1 + e^{-z}}$

tanh

ReLU

Leaky ReLU

One hidden layer
Neural Network

Why do you
need non-linear
activation functions?

# Activation function

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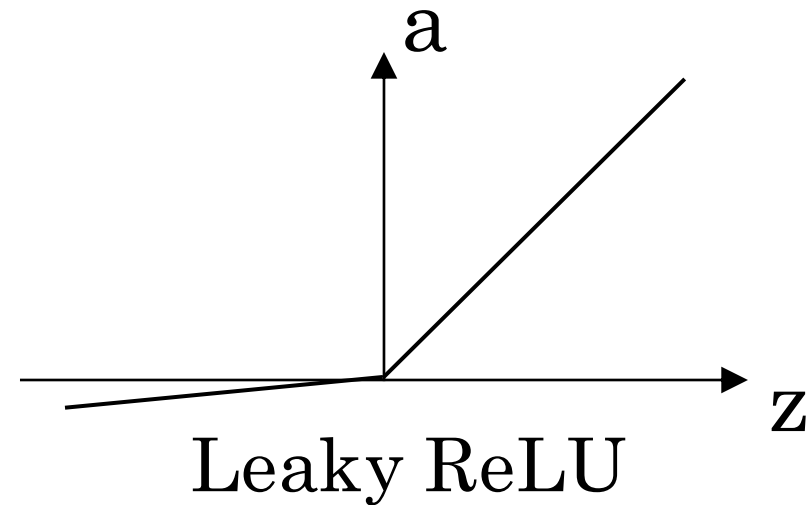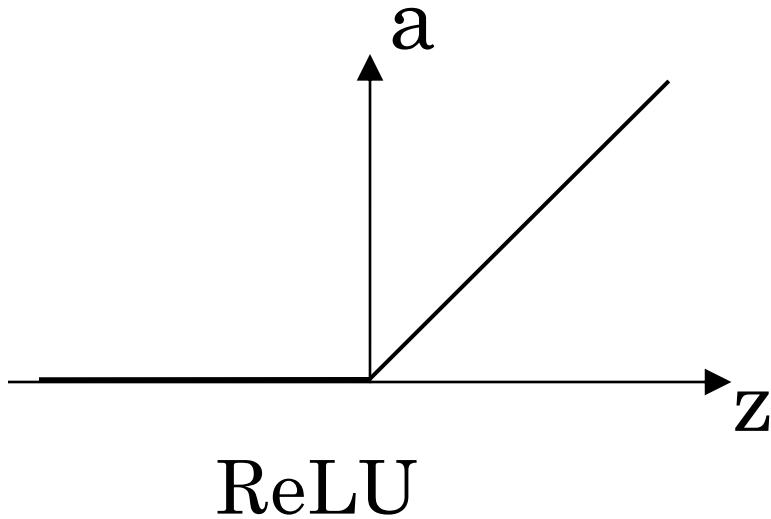
$$z^{[1]} = W^{[1]}x + b^{[1]}$$

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

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

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

# Activation function

Given x:

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

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

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

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

One hidden layer
Neural Network

Gradient descent for
neural networks

deeplearning.ai

# Gradient descent for neural networks

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

Cost function: $J\left(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}\right) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}, y)$

Repeat {

Compute predictions: $\left(\hat{y}^{(i)}, i = 1, \ldots m\right)$

$dW^{[1]} = \frac{\partial J}{\partial W^{[1]}} \quad , \quad db^{[1]} = \frac{\partial J}{\partial b^{[1]}} \quad , \ldots$

$W^{[1]} = W^{[1]} - \alpha \, dW^{[1]}$

$b^{[1]} = b^{[1]} - \alpha \, db^{[1]}$

$\ldots$ }

# Formulas for computing derivatives

## Forward propagation
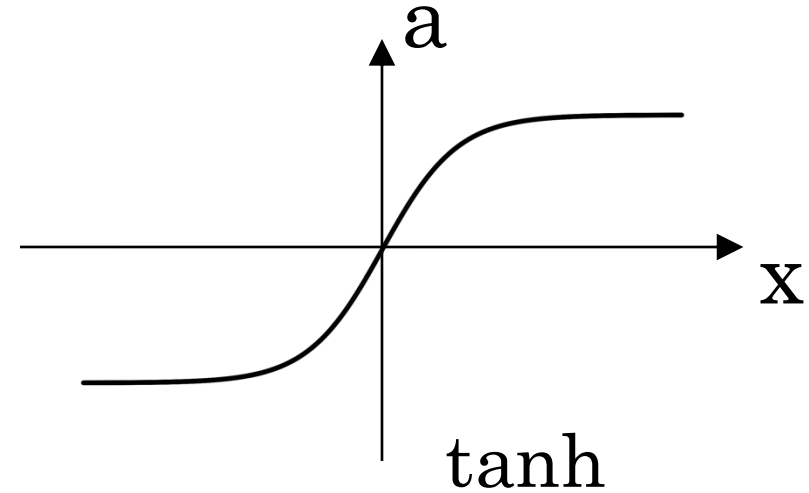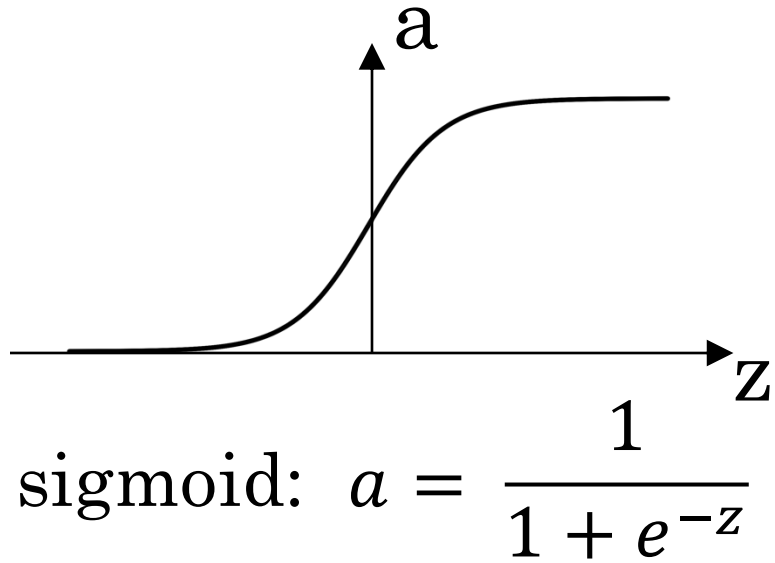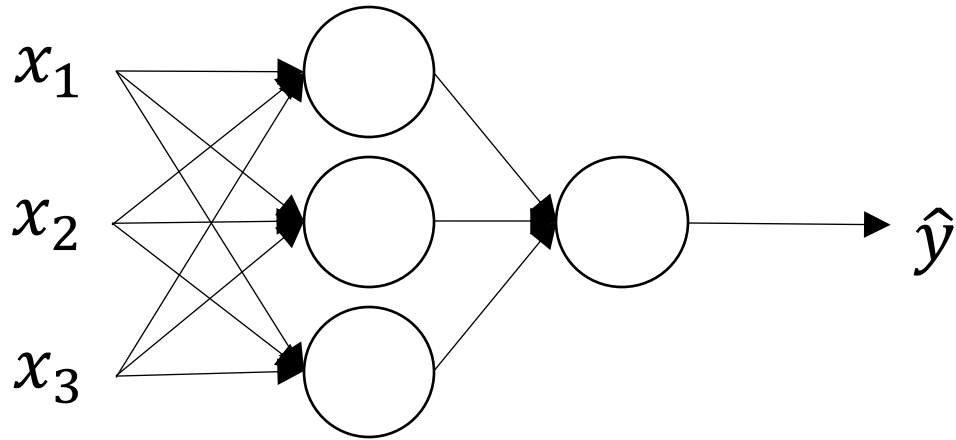
$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]})$$
$$= \sigma(Z^{[2]})$$

## Back propagation

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

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

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

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

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

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

One hidden layer
Neural Network

Backpropagation
intuition

deeplearning.ai

# Computing gradients

Logistic regression

$x$

$w$ $\rightarrow$ $z = w^T x + b$ $\rightarrow$ $a = \sigma(z)$ $\rightarrow$ $\mathcal{L}(a, y)$

$b$

$$da = \frac{d}{da}\mathcal{L}(a, y) = \frac{d}{da}\left(-y\log(a) - (1-y)\log(1-a)\right)$$
$$= -\frac{y}{a} + \frac{1-y}{1-a}$$

$$dz = da \cdot g'(z)$$

# Computing gradients

Logistic regression

$$x$$
$$w \longrightarrow \boxed{z = w^T x + b} \longrightarrow \boxed{a = \sigma(z)} \longrightarrow \boxed{\mathcal{L}(a, y)}$$
$$b$$

$$da = \frac{d}{da}\mathcal{L}(a, y) = \frac{d}{da}\left(-y\log(a) - (1 - y)\log(1 - a)\right)$$

$$= -\frac{y}{a} + \frac{1-y}{1-a}$$

$$dz = da \cdot g'(z)$$

to do: gradient of $g(z) = \frac{1}{1+e^{-z}}$

# Computing gradients

Logistic regression

$$x$$
$$w$$
$$b$$

$$z = w^T x + b$$ → $$a = \sigma(z)$$ → $$\mathcal{L}(a, y)$$

$$da = \frac{d}{da} \mathcal{L}(a, y) = \frac{d}{da} \left( -y \log(a) - (1-y) \log(1-a) \right)$$

$$= -\frac{y}{a} + \frac{1-y}{1-a}$$

$$dz = da \cdot g'(z)$$

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

Andrew Ng

# Computing gradients

Logistic regression

$$x$$
$$w$$
$$b$$

$$z = w^T x + b \quad\longrightarrow\quad a = \sigma(z) \quad\longrightarrow\quad \mathcal{L}(a, y)$$

$$da = \frac{d}{da}\mathcal{L}(a, y) = \frac{d}{da}\left(-y\log(a) - (1-y)\log(1-a)\right)$$

$$= -\frac{y}{a} + \frac{1-y}{1-a}$$

$$dz = da \cdot g'(z) = a - y$$
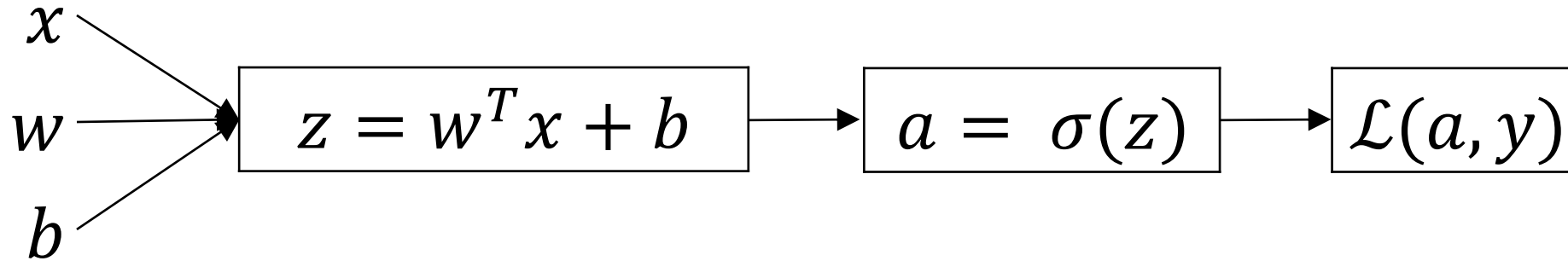
# Computing gradients

Logistic regression

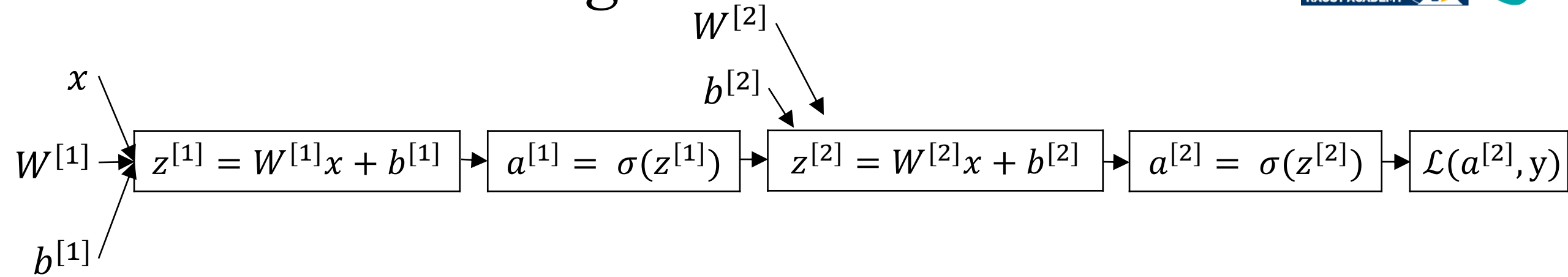$$x \searrow$$
$$w \rightarrow \boxed{z = w^T x + b} \rightarrow \boxed{a = \sigma(z)} \rightarrow \boxed{\mathcal{L}(a, y)}$$
$$b \nearrow$$

$$da = \frac{d}{da}\mathcal{L}(a, y) = \frac{d}{da}\left(-y\log(a) - (1 - y)\log(1 - a)\right)$$

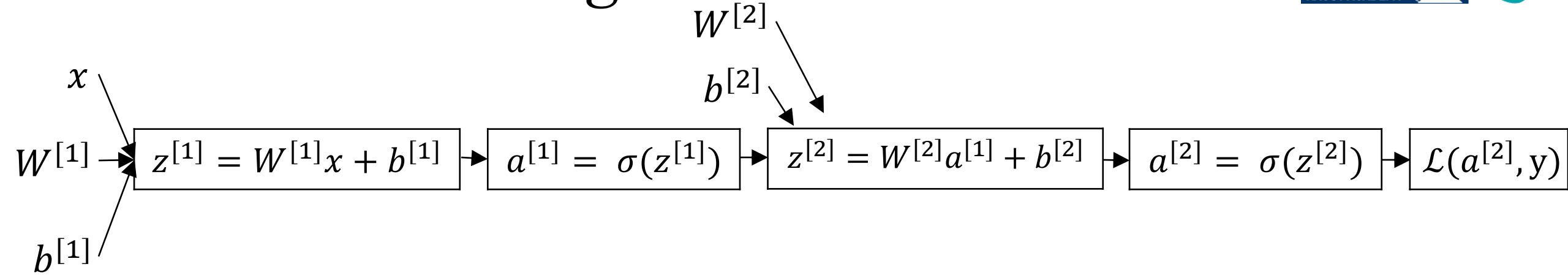$$= -\frac{y}{a} + \frac{1 - y}{1 - a}$$

$$dz = da \cdot g'(z) = a - y$$

$$dw = dz \cdot x$$

$$db = dz$$

Andrew Ng

# Neural network gradients



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

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

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

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

$$\mathcal{L}(a^{[2]}, y)$$

# Neural network gradients

$$W^{[2]}$$

$$b^{[2]}$$

$$x$$

$$W^{[1]} \rightarrow$$

$$b^{[1]}$$

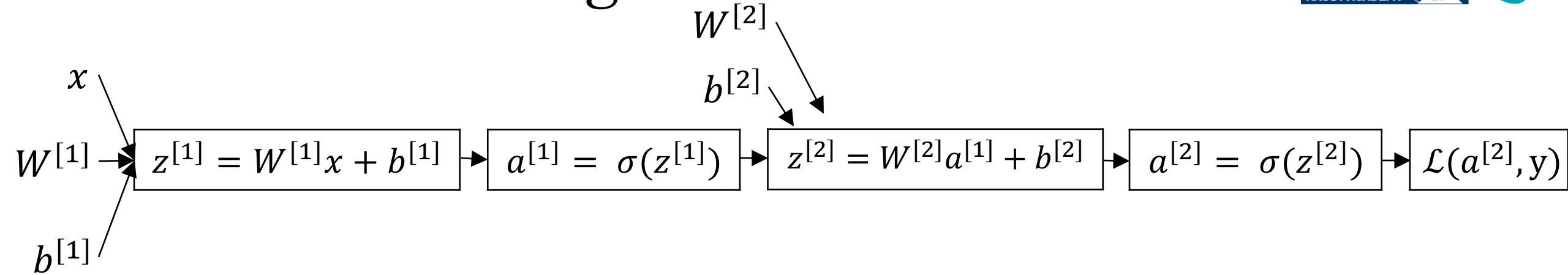| $z^{[1]} = W^{[1]}x + b^{[1]}$ | $a^{[1]} = \sigma(z^{[1]})$ | $z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$ | $a^{[2]} = \sigma(z^{[2]})$ | $\mathcal{L}(a^{[2]}, y)$ |

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

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

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

Andrew Ng

# Neural network gradients

$$z^{[1]} = W^{[1]}x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow \mathcal{L}(a^{[2]}, y)$$

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

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

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

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

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

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

Andrew Ng

# Summary of gradient descent

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

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

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

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

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

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

# Summary of gradient descent

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

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

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

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

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

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

Vectorization implementation

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

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$
$$A^{[1]} = \sigma(Z^{[1]})$$
$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$
$$A^{[2]} = \sigma(Z^{[2]})$$

Andrew Ng

# Summary of gradient descent

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

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

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

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

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

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

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

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

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

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

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

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

Andrew Ng

# Scale drives deep learning progress

# Scale drives deep learning progress

- Data

- Computation

- Algorithms



Idea

Code

Experiment

# What happens if you initialize weights to zero?



Initial weights = 0 → symmetry → similar updates

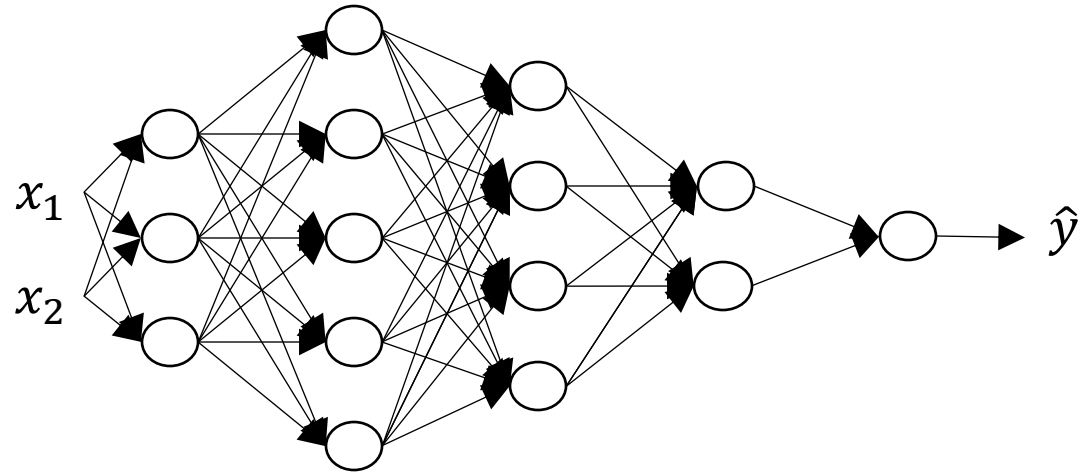# Random initialization



small values for $W^{[1]}$ and $W^{[2]}$

# Deep Neural Networks

deeplearning.ai

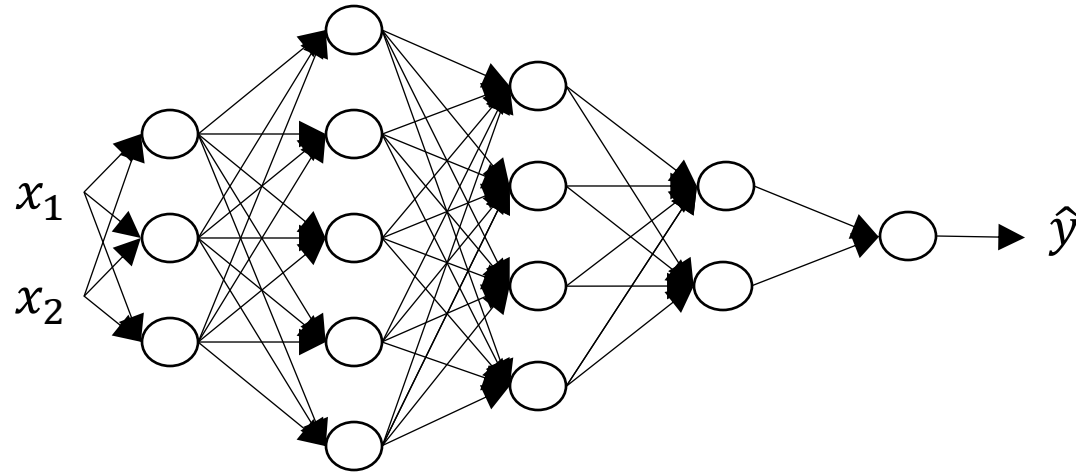## Getting your matrix dimensions right

# Parameters $W^{[l]}$ and $b^{[l]}$



$$W^{[l]} : (n^{[l]}, n^{[l-1]})$$

$$b^{[l]} : (n^{[l]}, 1)$$

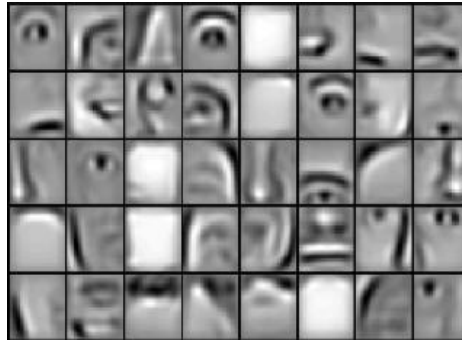# Vectorized implementation
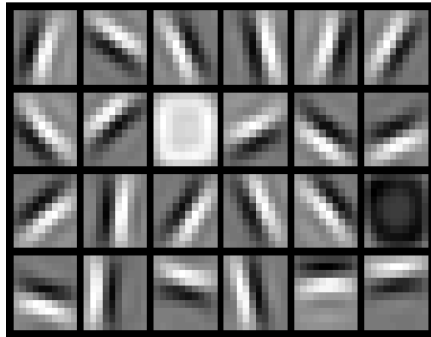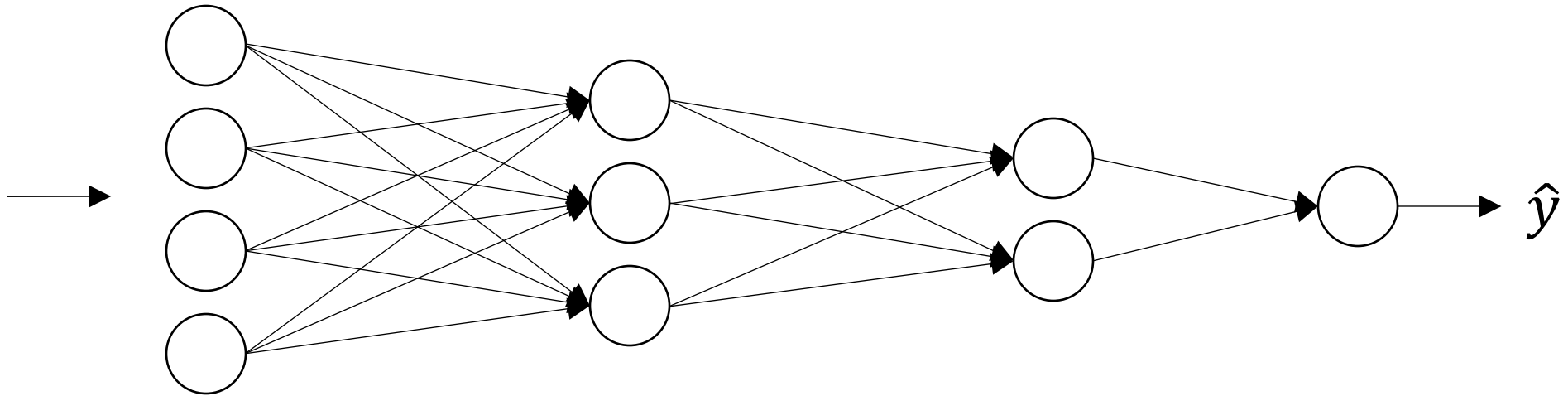


$$Z^{[l]}, A^{[l]} : (n^{[l]}, m)$$

# Deep Neural Networks

deeplearning.ai

## Why deep representations?

# Intuition about deep representation

# 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.
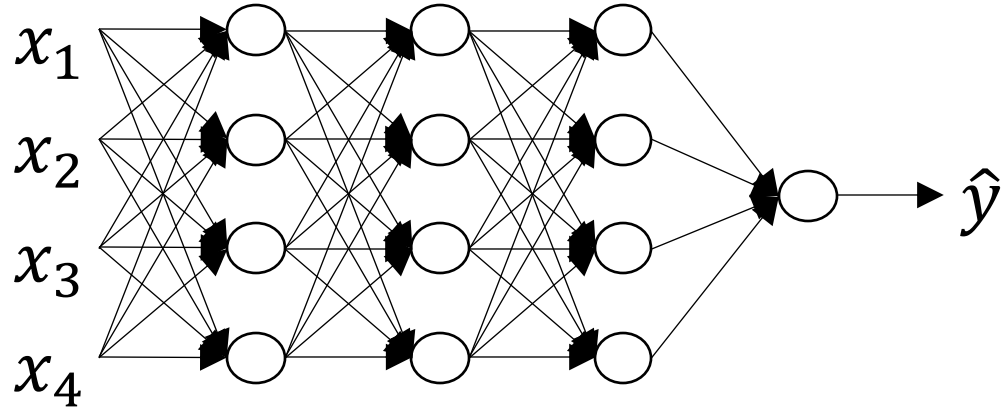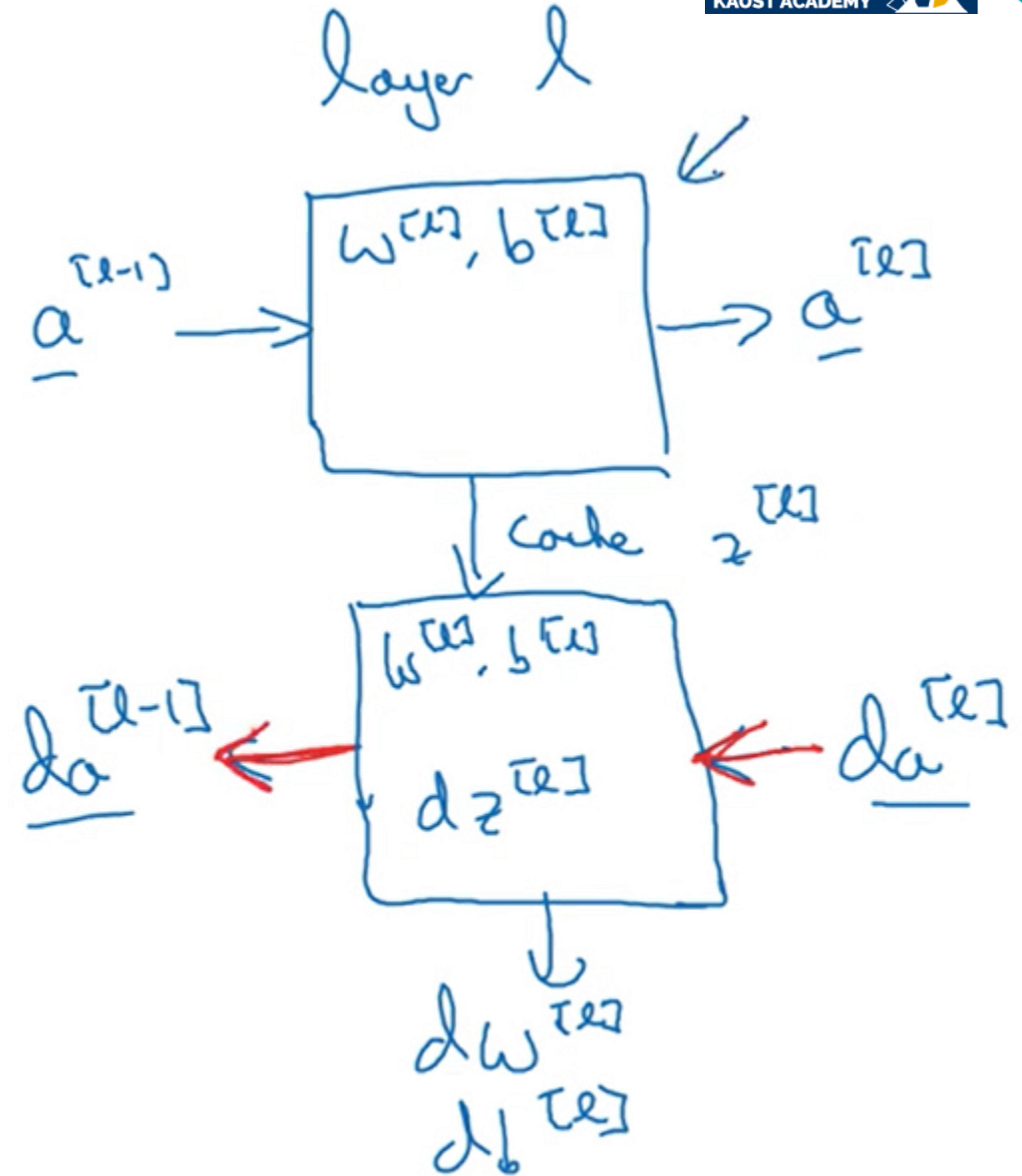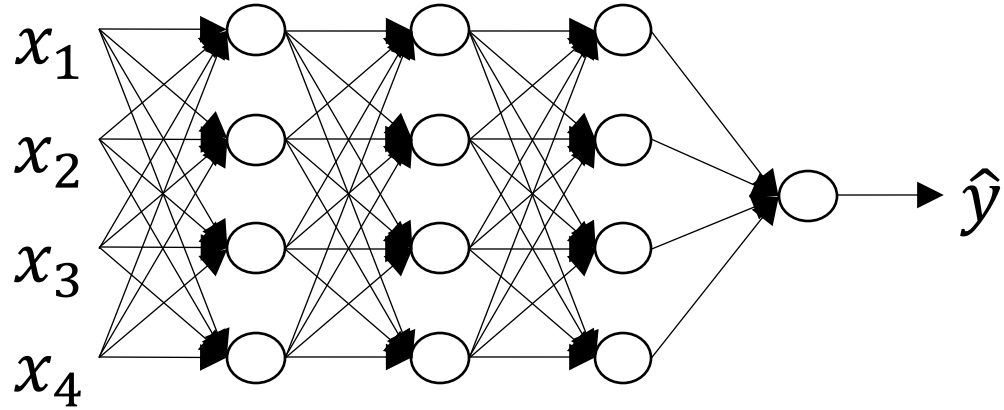
Example: xor

# Deep Neural Networks

deeplearning.ai

## Building blocks of deep neural networks

# Forward and backward functions

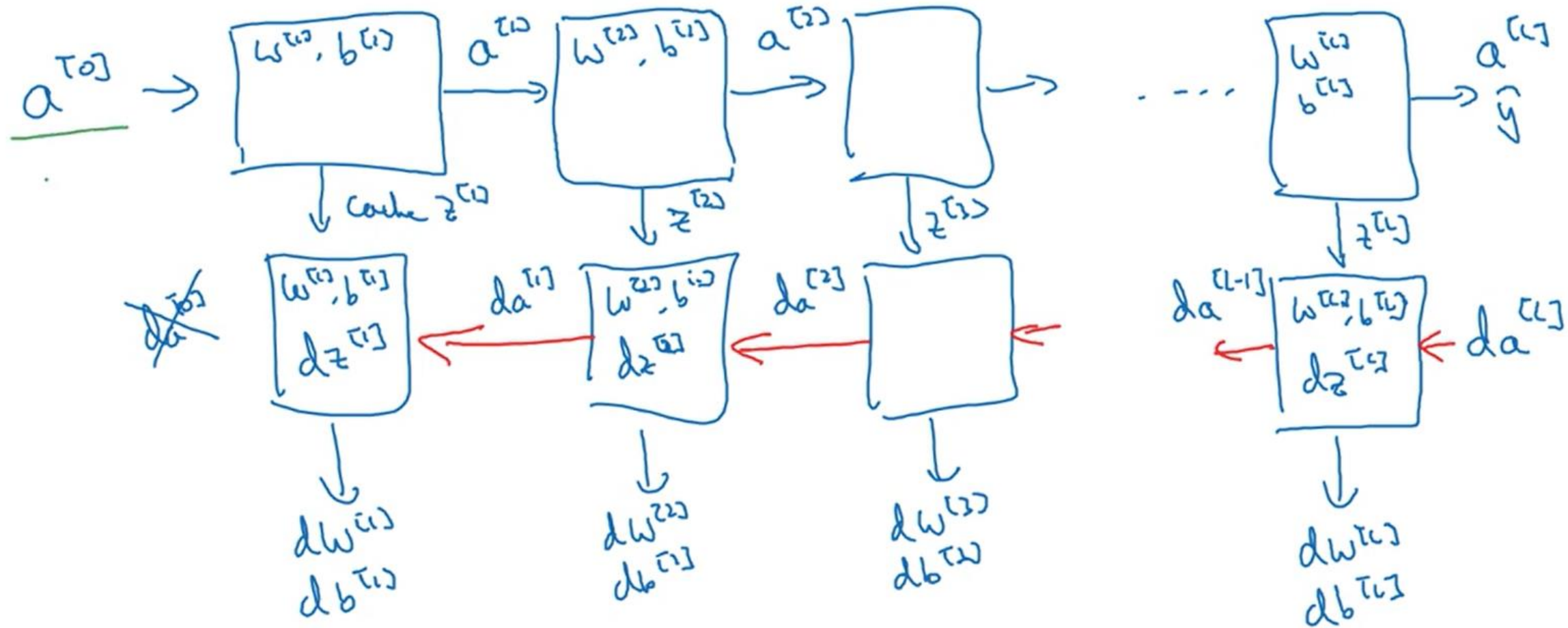# Forward and backward functions



Andrew Ng

# Forward and backward functions

deeplearning.ai

# Deep Neural Networks

## Forward and backward propagation

# Forward propagation for layer *l*

Input $a^{[l-1]}$

Output $a^{[l]}$, cache $(z^{[l]})$

# Forward propagation for layer *l*

Input $a^{[l-1]}$

Output $a^{[l]}$, cache $(z^{[l]})$

$$Z^{[l]} = W^{[l]}A^{[l-1]} + b^{[l]}$$
$$A^{[l]} = g^{[l]}(Z^{[l]})$$

# Backward propagation for layer *l*

Input $da^{[l]}$

Output $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]\prime}(z^{[1]})$$

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

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

Previously

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

$$dW^{[l]} = dz^{[l]} a^{[l-1]}$$

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

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