



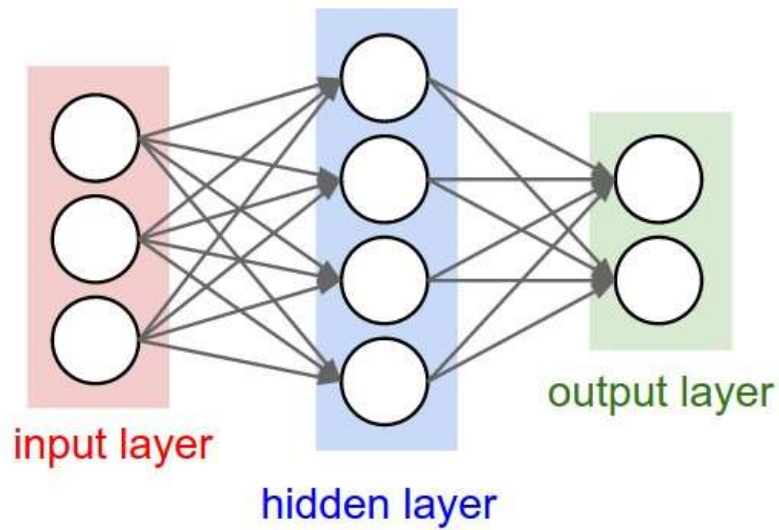
deeplearning.ai

# Introduction to Deep Learning

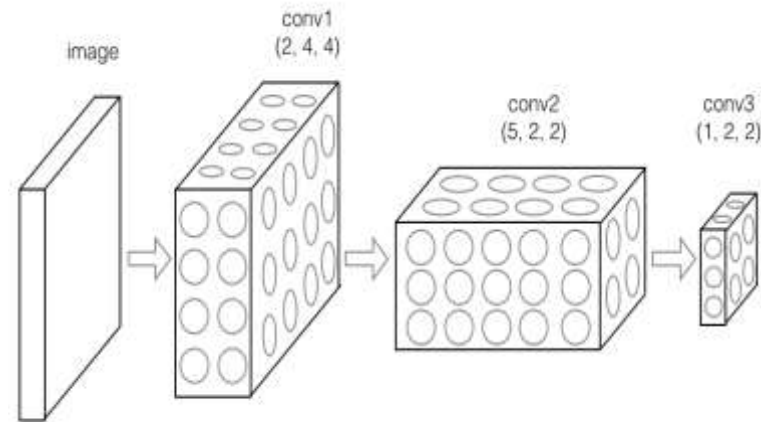
---

## Supervised Learning with Neural Networks

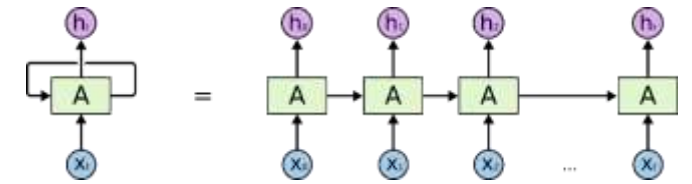
# Neural Network examples



Standard NN



Convolutional NN

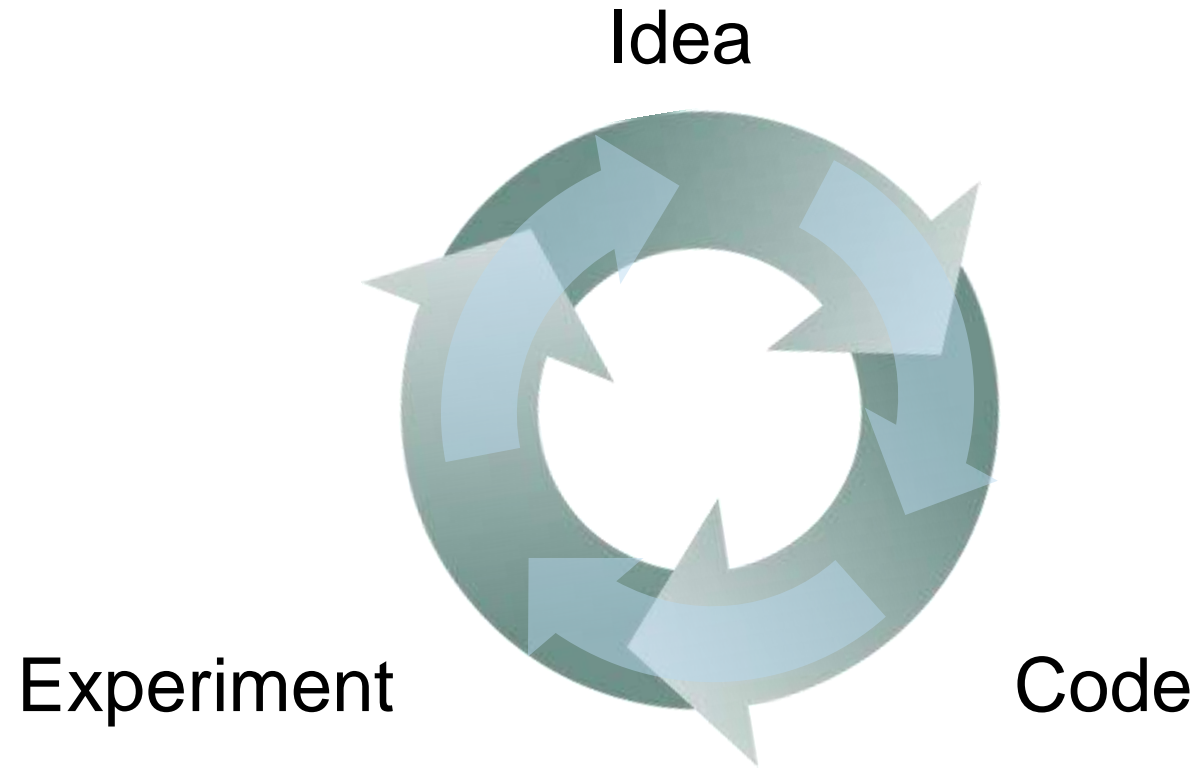


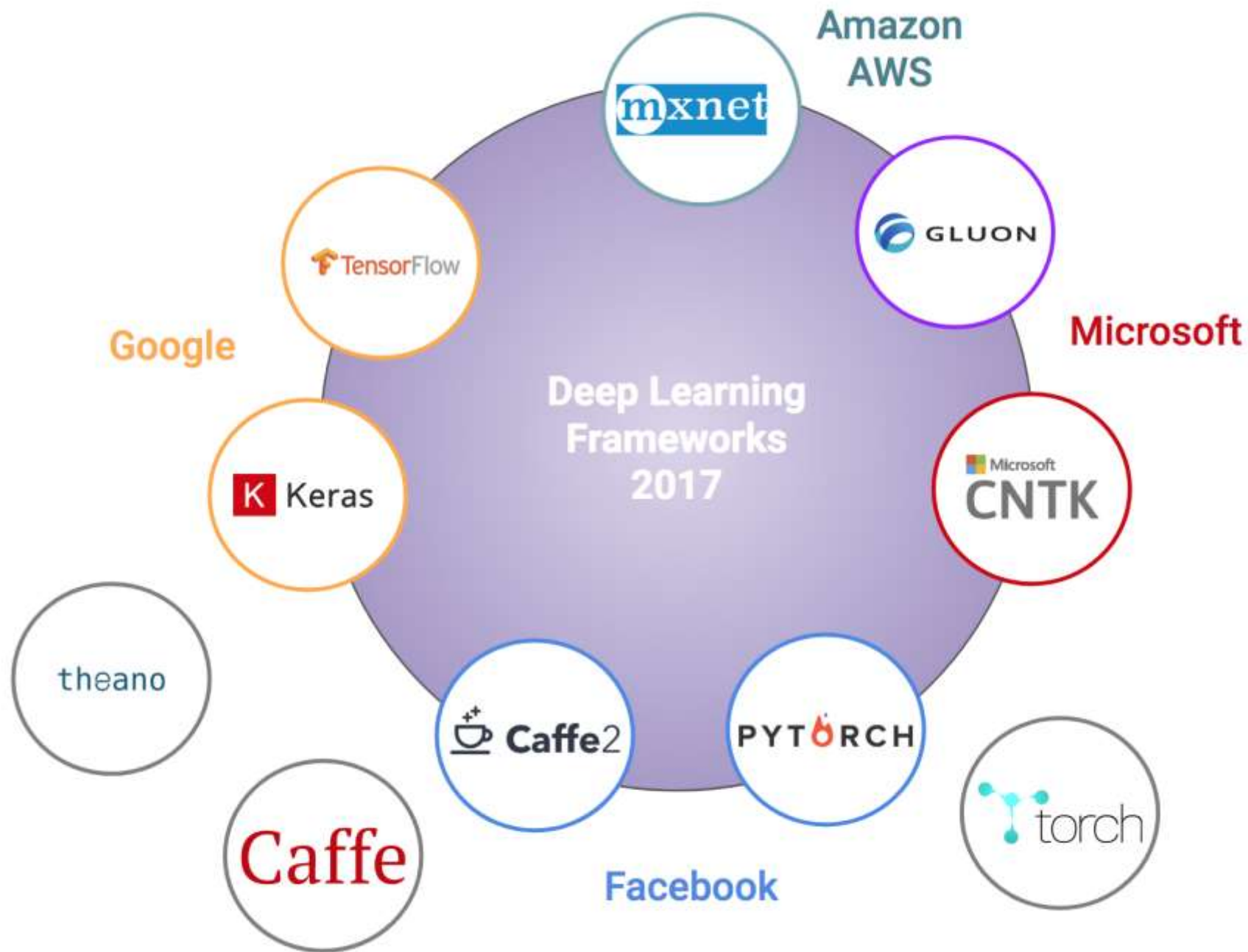
Recurrent NN

# Why is Deep Learning taking off?

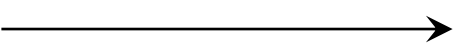
- More data
- More computational resources
- New algorithms

# Deep learning is highly iterative process





# Binary Classification



1 (cat) vs 0 (non cat)

		Blue			
Green		255	134	93	22
Red		255	134	202	22
	255	231	42	22	4
	123	94	83	2	192
	34	44	187	92	34
	34	76	232	124	94
	67	83	194	202	

# Notation

# Logistic Regression



# Logistic Regression cost function

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

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

Loss (error) function:

# Logistic regression cost function

$$\begin{aligned} \rightarrow & \text{If } y = 1: p(y|x) = \hat{y} \\ \rightarrow & \text{If } y = 0: p(y|x) = 1 - \hat{y} \end{aligned} \quad \left. \vphantom{\begin{aligned} \rightarrow & \text{If } y = 1: p(y|x) = \hat{y} \\ \rightarrow & \text{If } y = 0: p(y|x) = 1 - \hat{y} \end{aligned}} \right\} p(y|x)$$

$$p(y|x) = \hat{y}^y (1 - \hat{y})^{(1-y)}$$

$$\text{If } y=1: p(y|x) = \hat{y} \underbrace{(1-\hat{y})^0}_{=1}$$

$$\text{If } y=0: p(y|x) = \hat{y}^0 \underbrace{(1-\hat{y})^{(1-0)}}_{=1} = 1 \times (1-\hat{y}) = 1 - \hat{y}$$

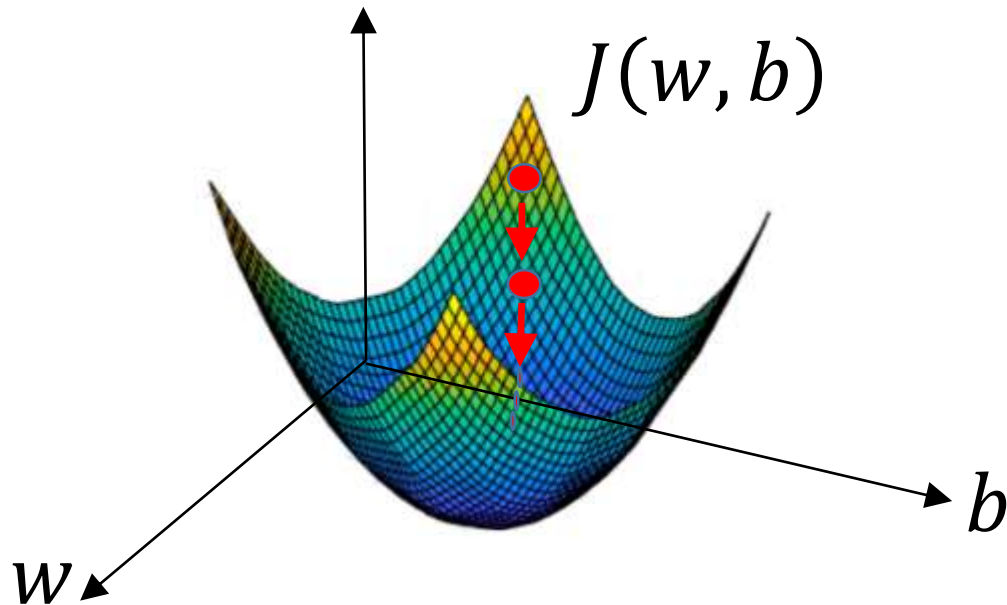
$$\begin{aligned} \log p(y|x) &= \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ &= -\ell(\hat{y}, y) \end{aligned}$$

# Gradient Descent

Recap:  $\hat{y} = \sigma(w^T x + b)$ ,  $\sigma(z) = \frac{1}{1+e^{-z}}$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

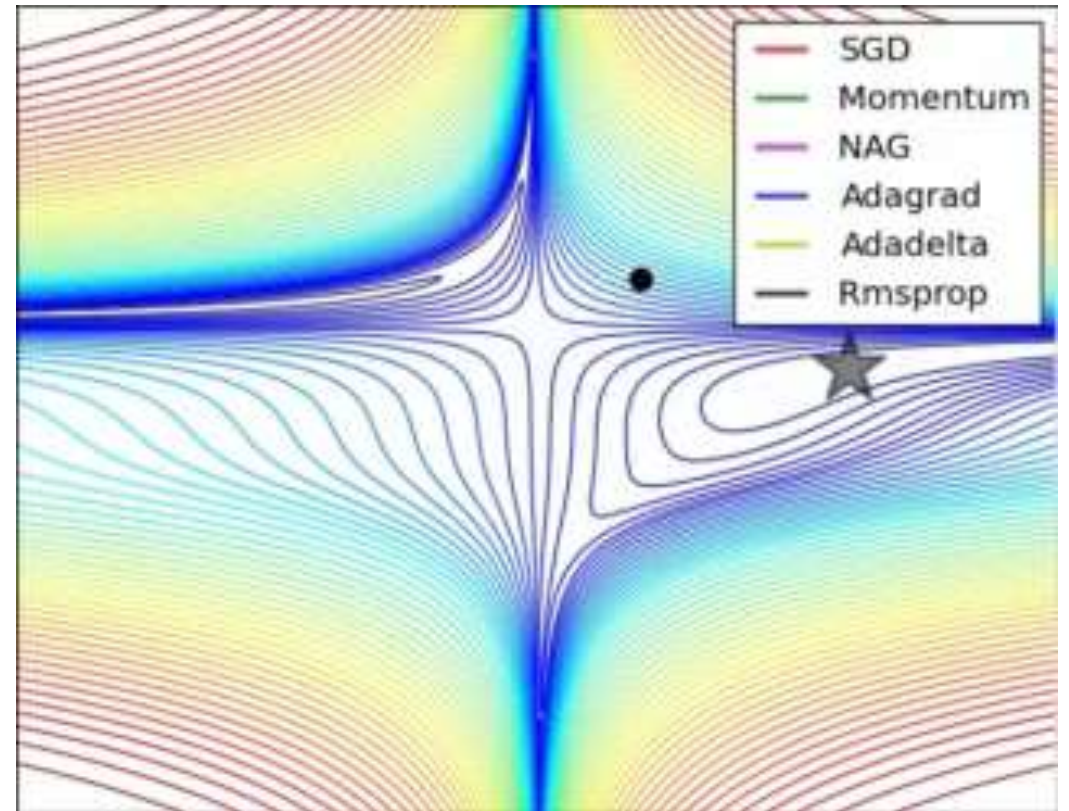
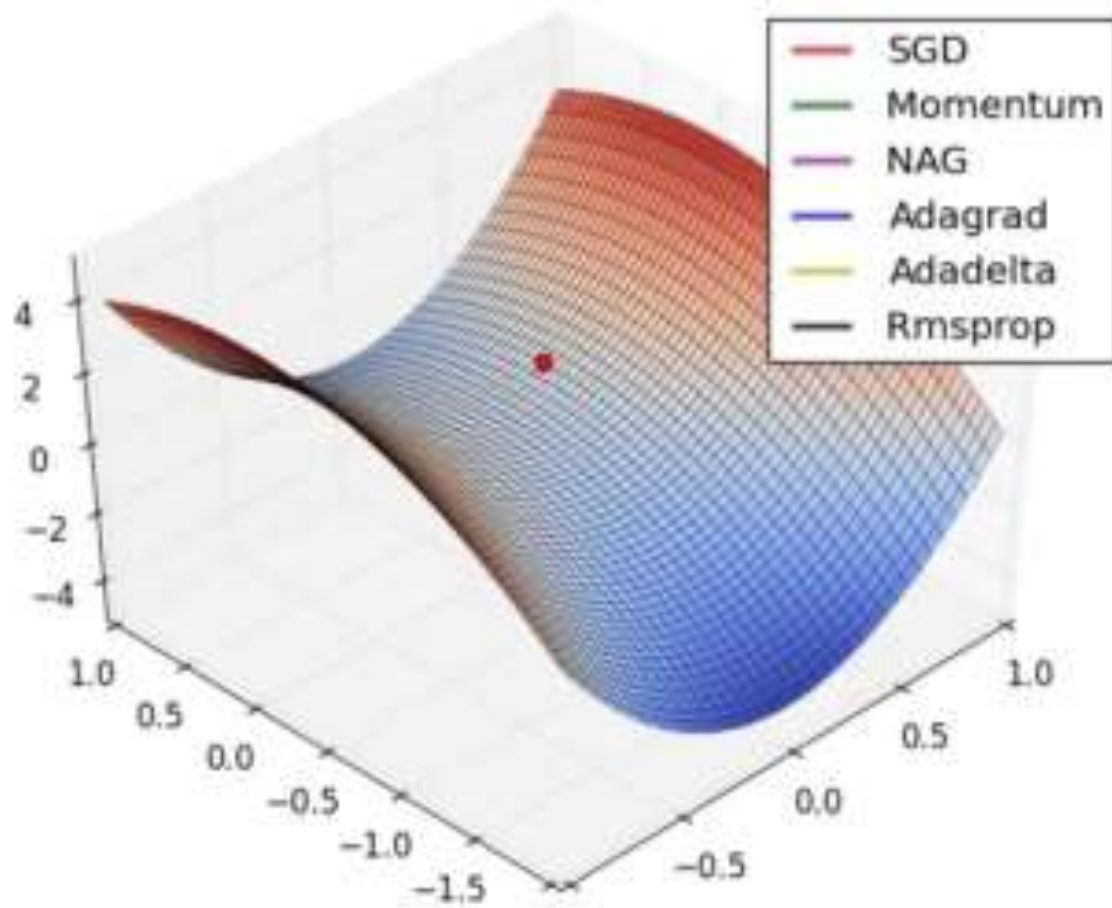
Want to find  $w, b$  that minimize  $J(w, b)$



# Gradient Descent



# Optimization methods





deeplearning.ai

# Basics of Neural Network Programming

---

## Computation Graph

# Derivatives

$$f(x, y) = x + y \rightarrow \frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial y} = 1$$

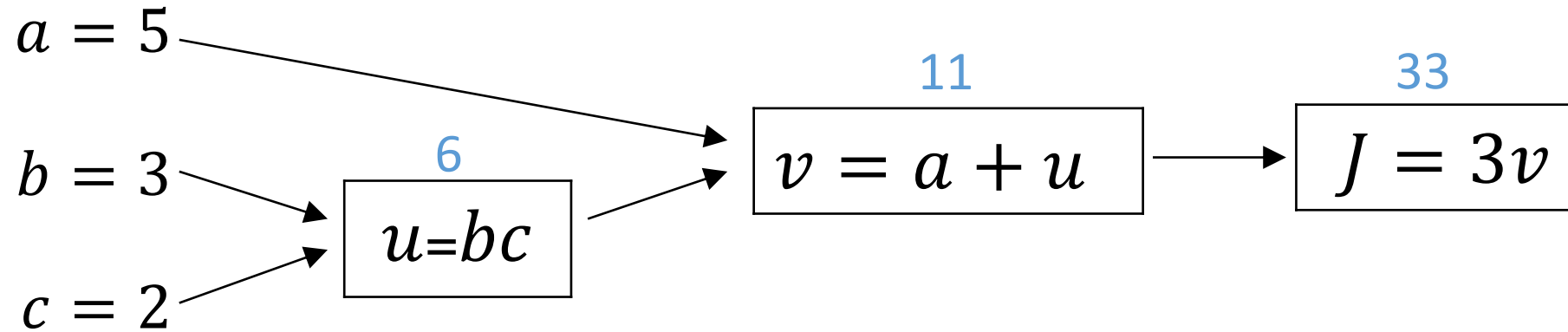
$$f(x, y) = x \cdot y \rightarrow \frac{\partial f}{\partial x} = y \quad \frac{\partial f}{\partial y} = x$$

$$f(x) = y(z(x)) \rightarrow \frac{df}{dx} = \frac{df}{dy} \cdot \frac{dy}{dz} \cdot \frac{dz}{dx}$$

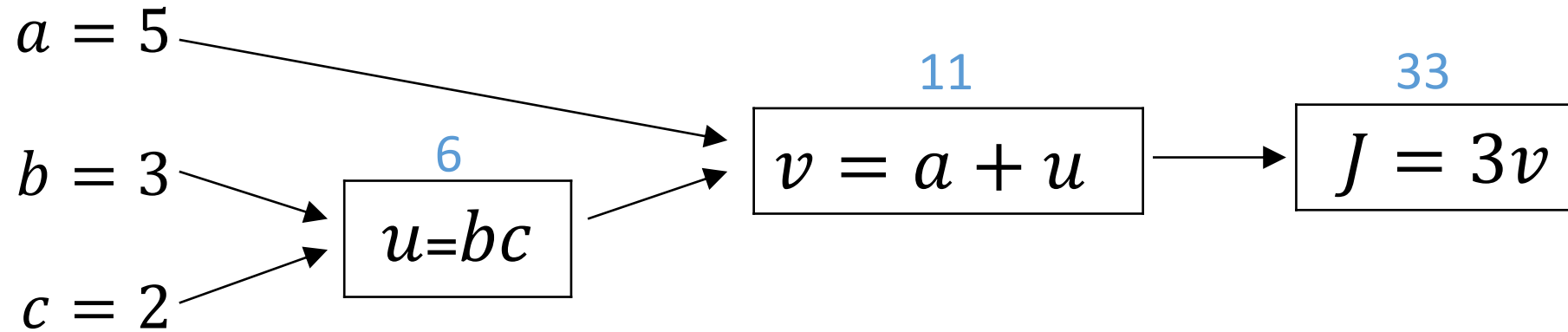
# Computation Graph



# Computing derivatives

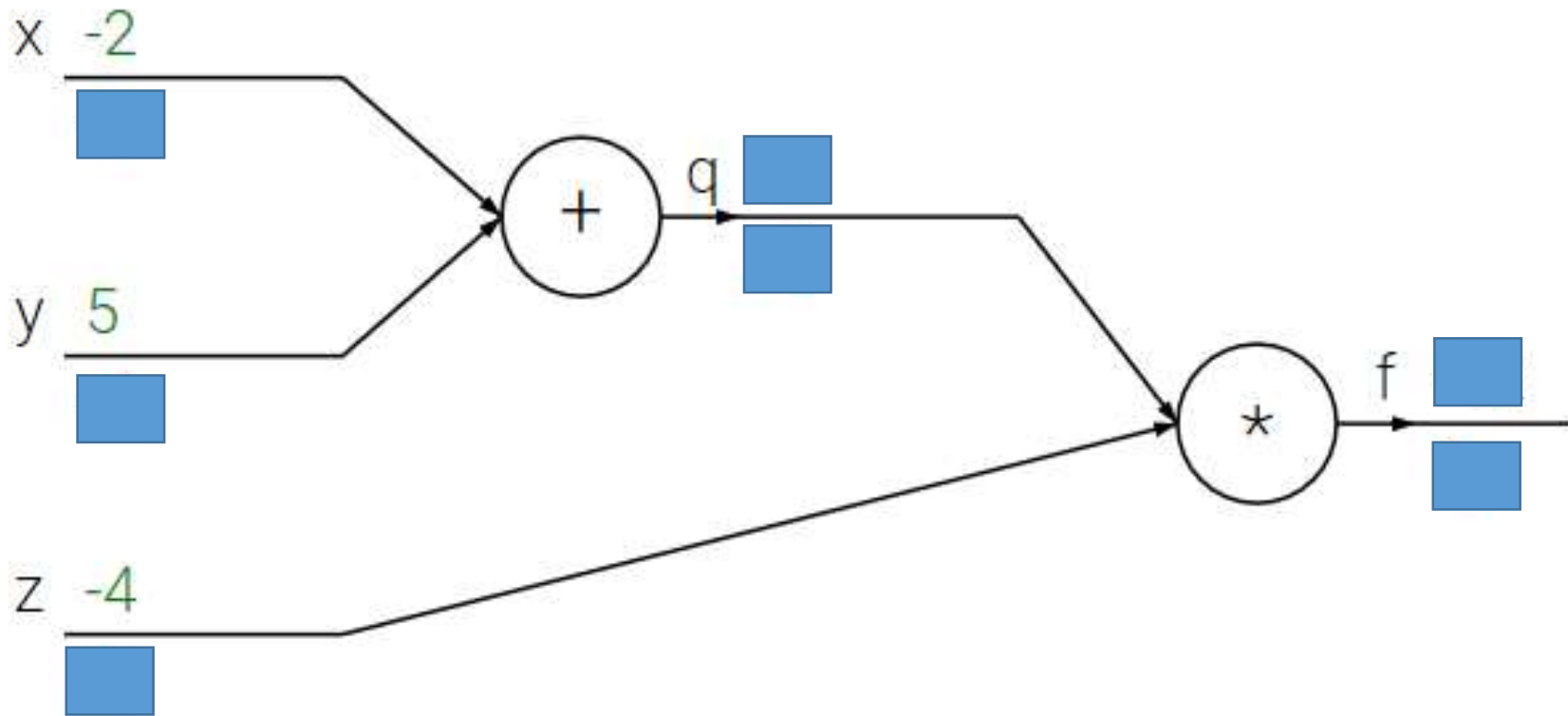


# Computing derivatives



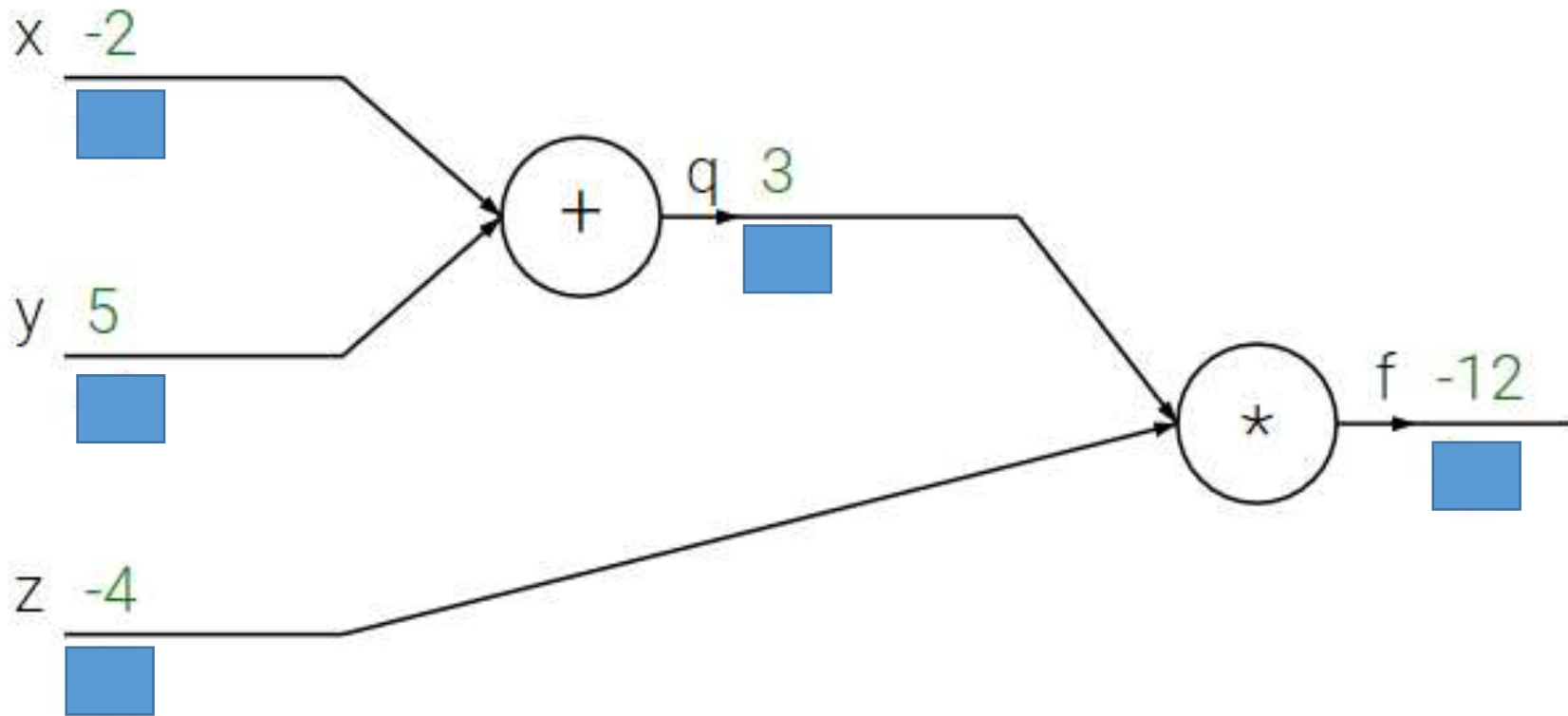
# Computational graph

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



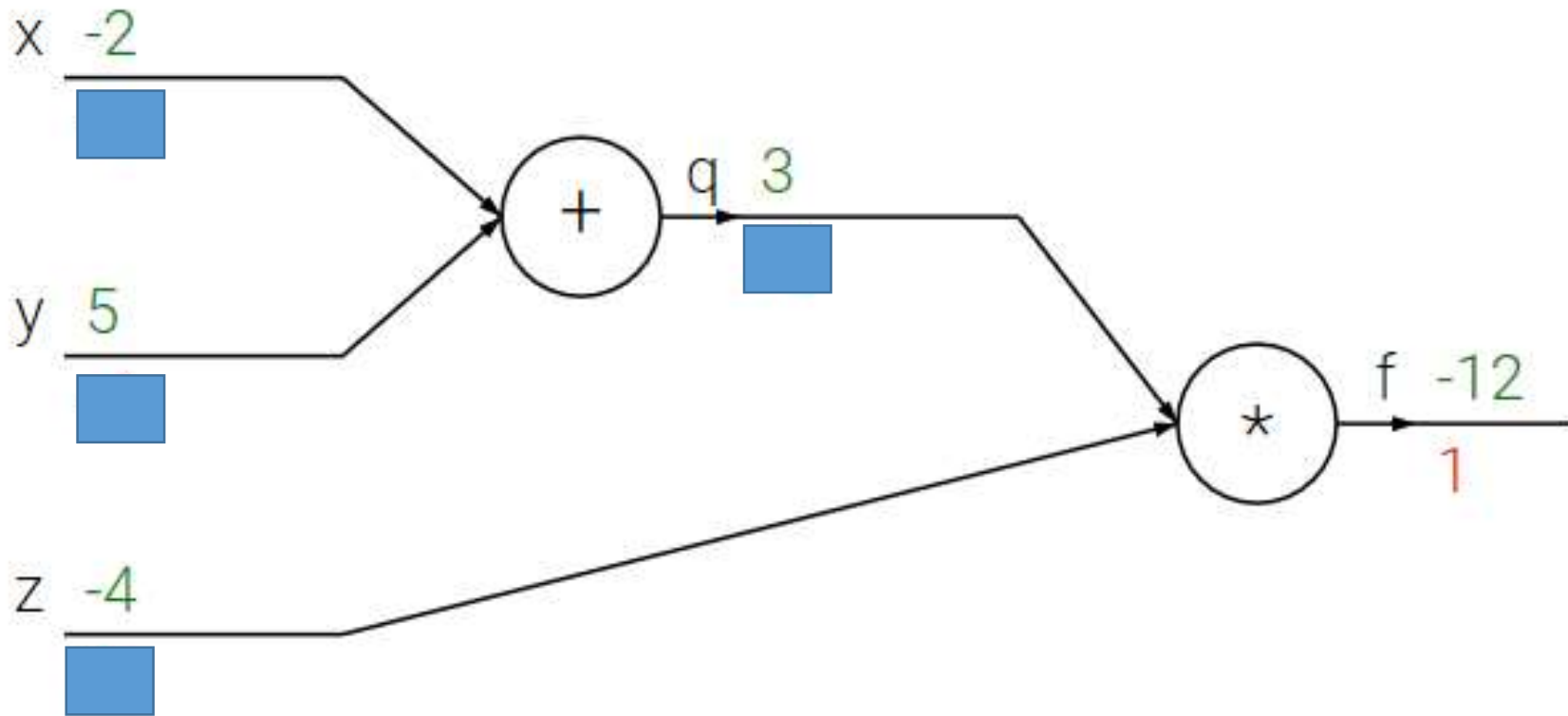
# Computational graph

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



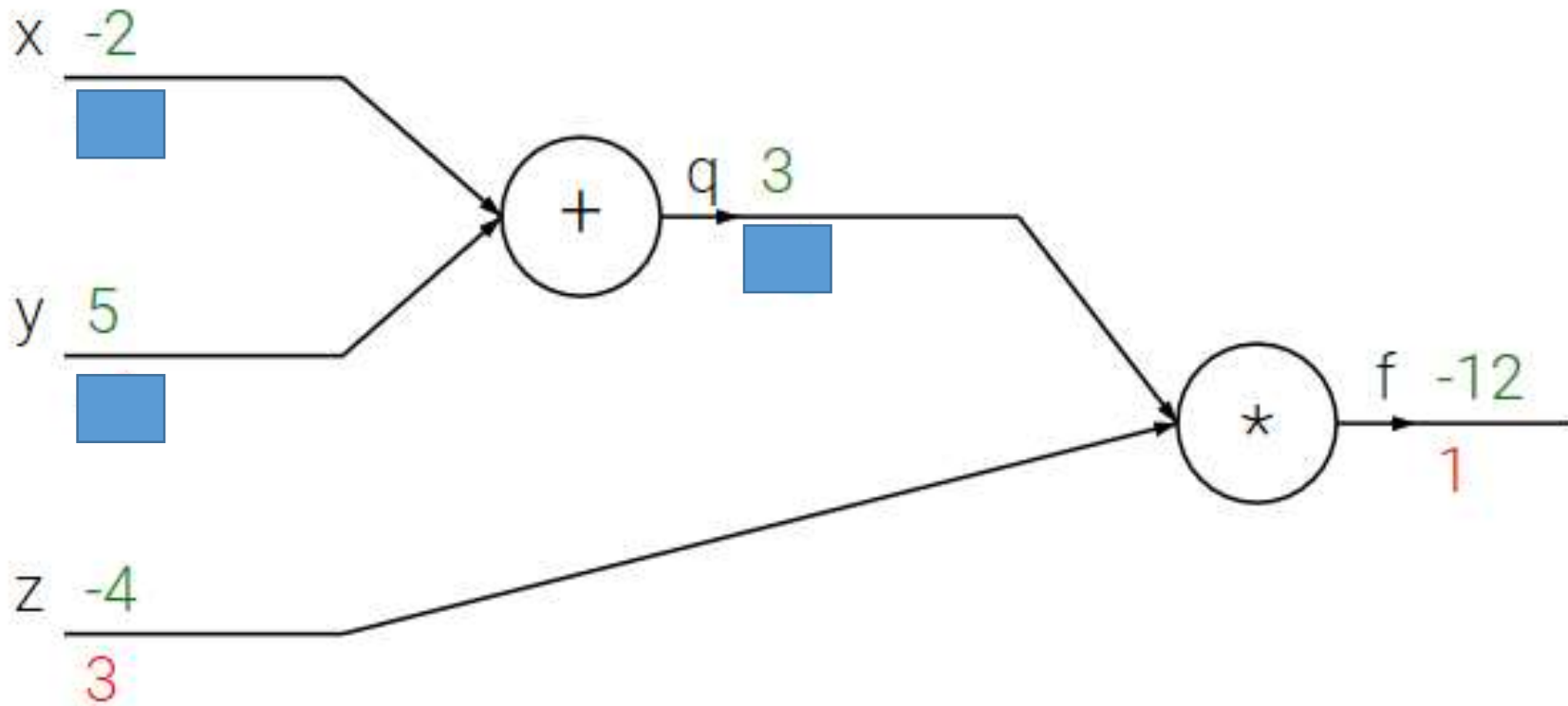
# Computational graph

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



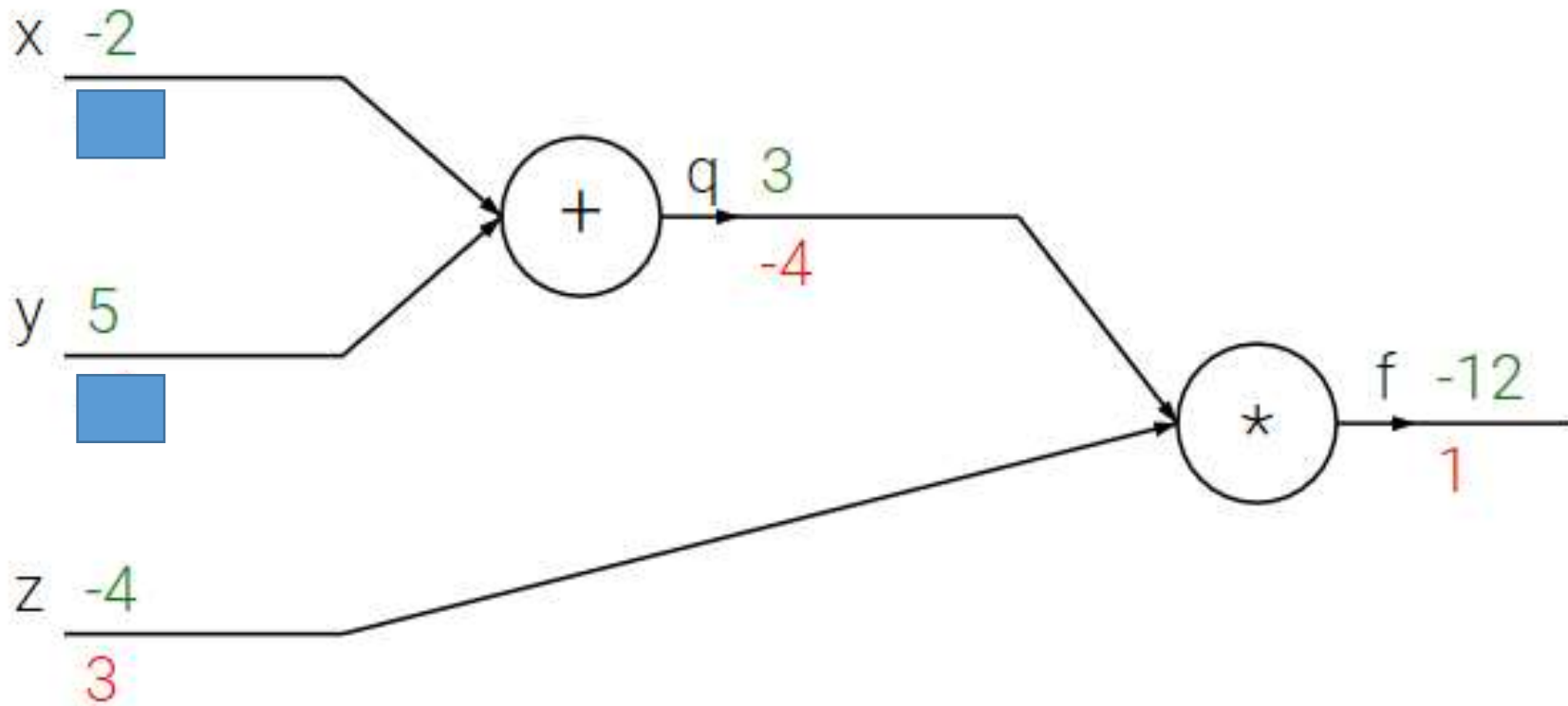
# Computational graph

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



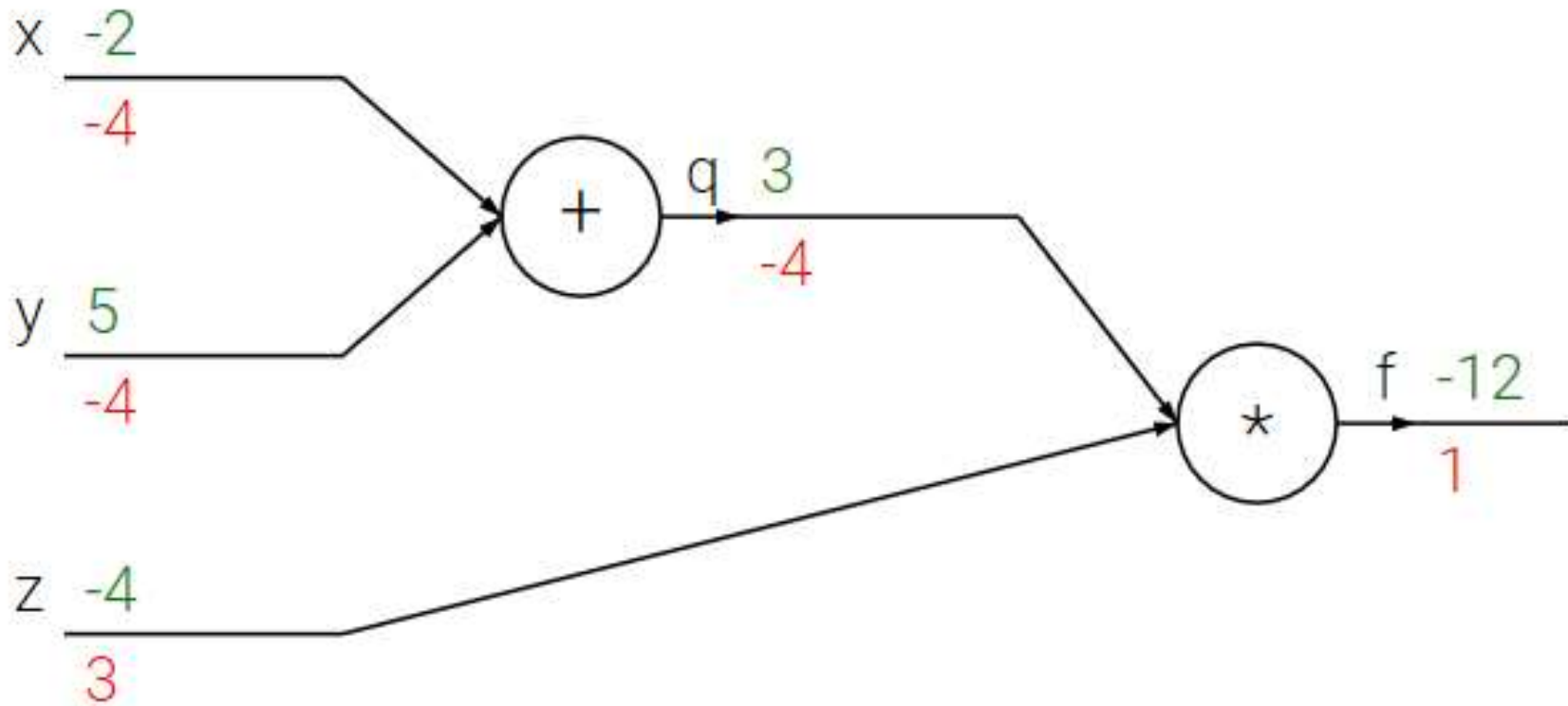
# Computational graph

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



# Computational graph

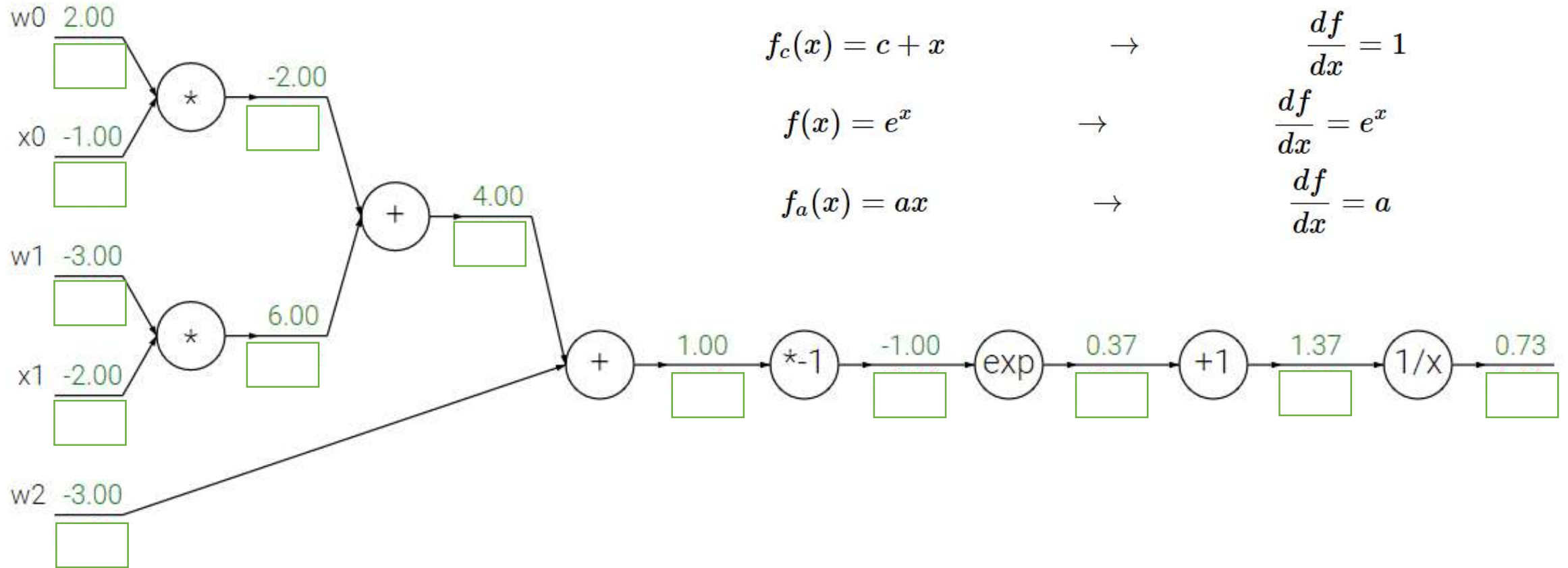
$$f(x, y, z) = (x + y) * z$$





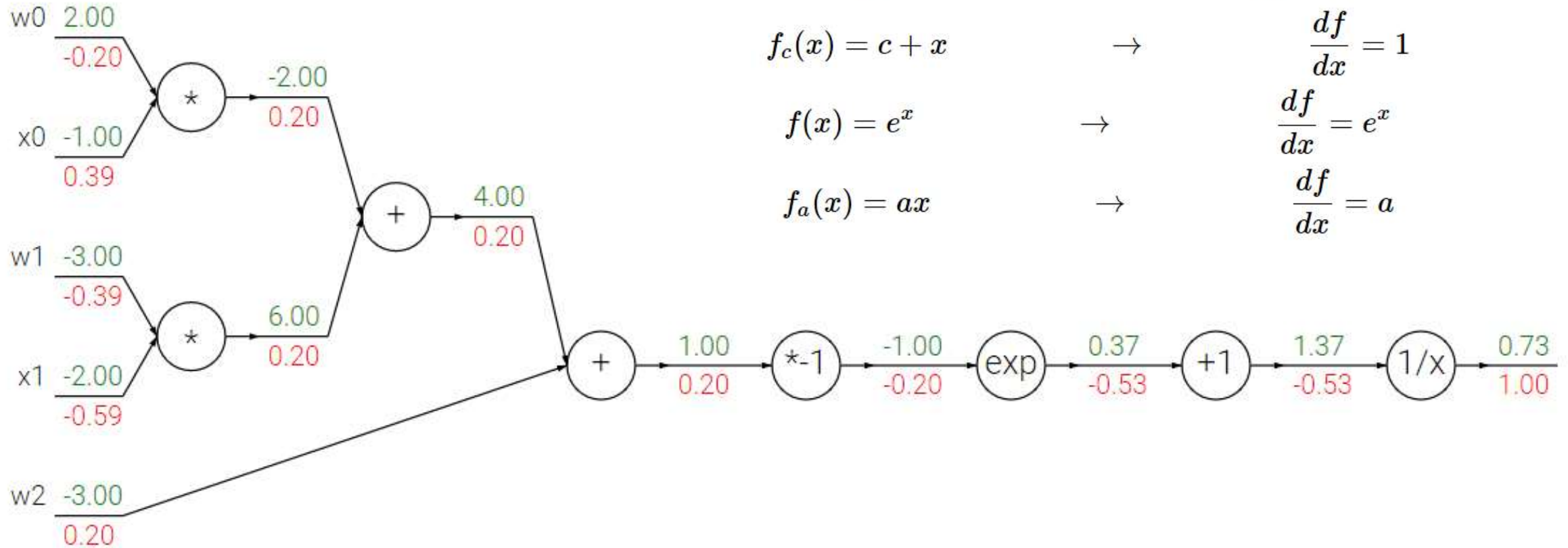
# Logistic regression

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



# Logistic regression

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$



$$f(x) = \frac{1}{x}$$

→

$$\frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x$$

→

$$\frac{df}{dx} = 1$$

$$f(x) = e^x$$

→

$$\frac{df}{dx} = e^x$$

$$f_a(x) = ax$$

→

$$\frac{df}{dx} = a$$



deeplearning.ai

# Basics of Neural Network Programming

---

Logistic Regression  
Gradient descent

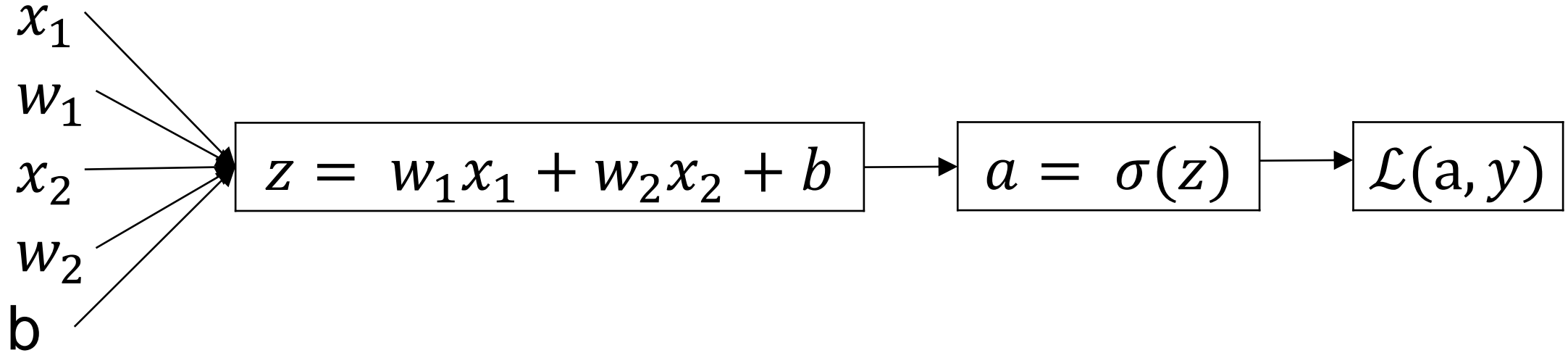
# Logistic regression recap

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

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

$$\mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$$

# Logistic regression derivatives



# Logistic regression on $m$ examples

# Logistic regression on $m$ examples



deeplearning.ai

# Basics of Neural Network Programming

---

## Vectorization



# What is vectorization?

# Neural network programming guideline

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.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$$

```
u = np.zeros( (n,1) )  
for i in range(n):  
    u[i]=math.exp(v[i])
```

# Logistic regression derivatives

$$J = 0, \quad dw_1 = 0, \quad dw_2 = 0, \quad db = 0$$

for  $i = 1$  to  $n$ :

$$z^{(i)} = w^T x^{(i)} + b$$

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

$$J += -[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

$$dz^{(i)} = a^{(i)}(1 - a^{(i)})$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J = J/m, \quad dw_1 = dw_1/m, \quad dw_2 = dw_2/m, \quad db = db/m$$



deeplearning.ai

# Basics of Neural Network Programming

---

## Vectorizing Logistic Regression

# Vectorizing Logistic Regression

$$z^{(1)} = w^T x^{(1)} + b$$

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

$$z^{(2)} = w^T x^{(2)} + b$$

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

$$z^{(3)} = w^T x^{(3)} + b$$

$$a^{(3)} = \sigma(z^{(3)})$$

# Vectorizing Logistic Regression

# Implementing Logistic Regression

$J = 0, \quad dw_1 = 0, \quad dw_2 = 0, \quad db = 0$

for  $i = 1$  to  $m$ :

$$z^{(i)} = w^T x^{(i)} + b$$

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

$$J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$J = J/m, \quad dw_1 = dw_1/m, \quad dw_2 = dw_2/m$

$db = db/m$





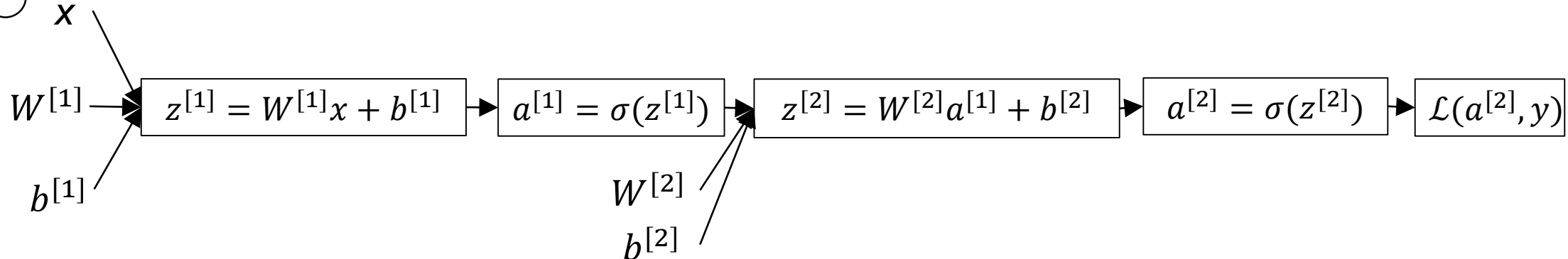
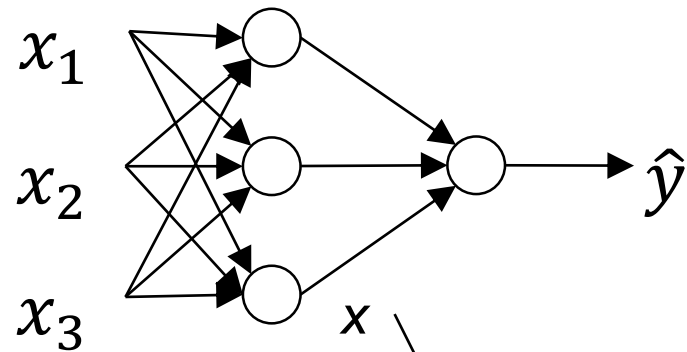
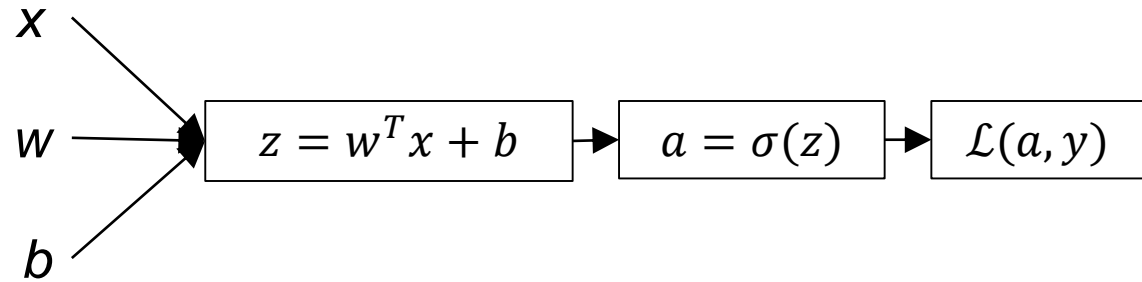
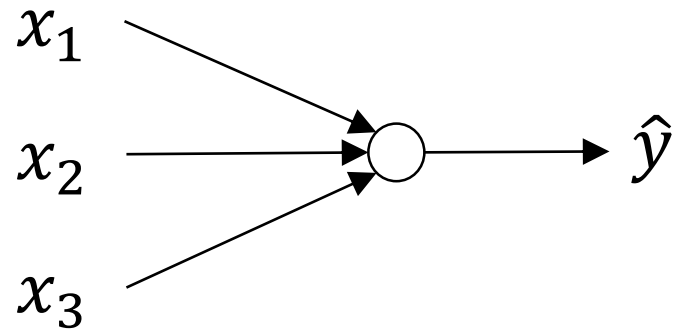
deeplearning.ai

One hidden layer  
Neural Network

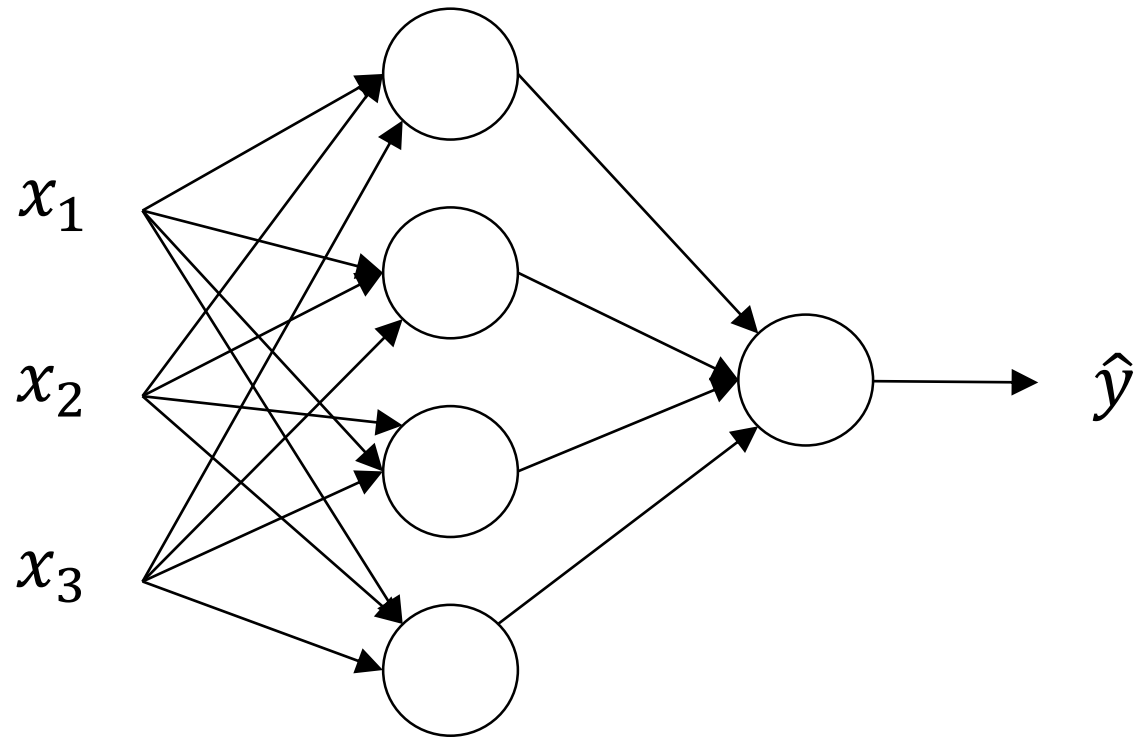
---

# Neural Networks Overview

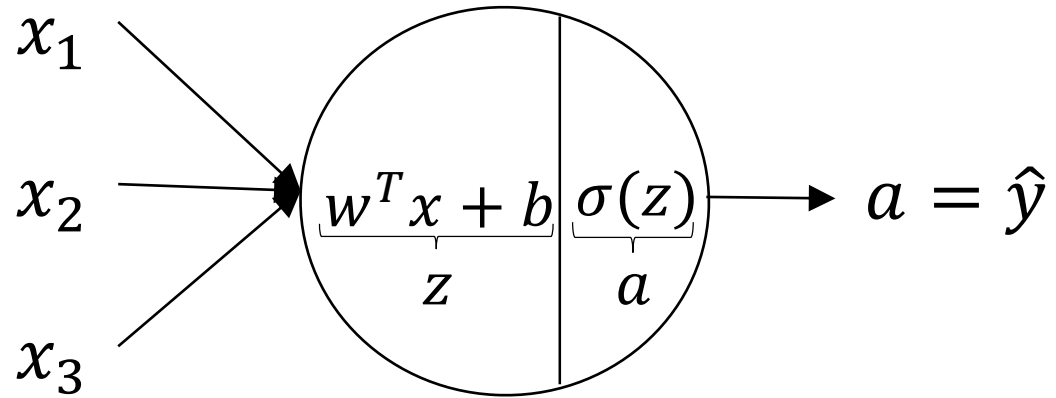
# What is a Neural Network?



# Neural Network Representation

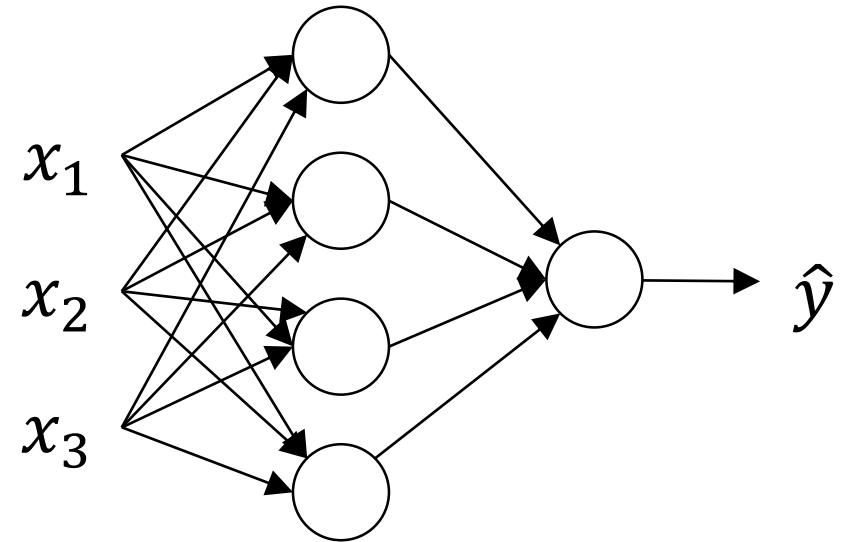


# Neural Network Representation

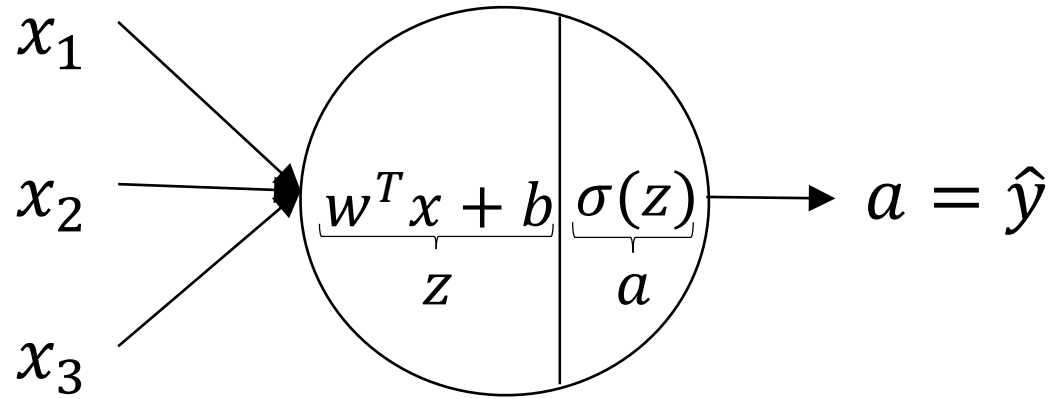


$$z = w^T x + b$$

$$a = \sigma(z)$$

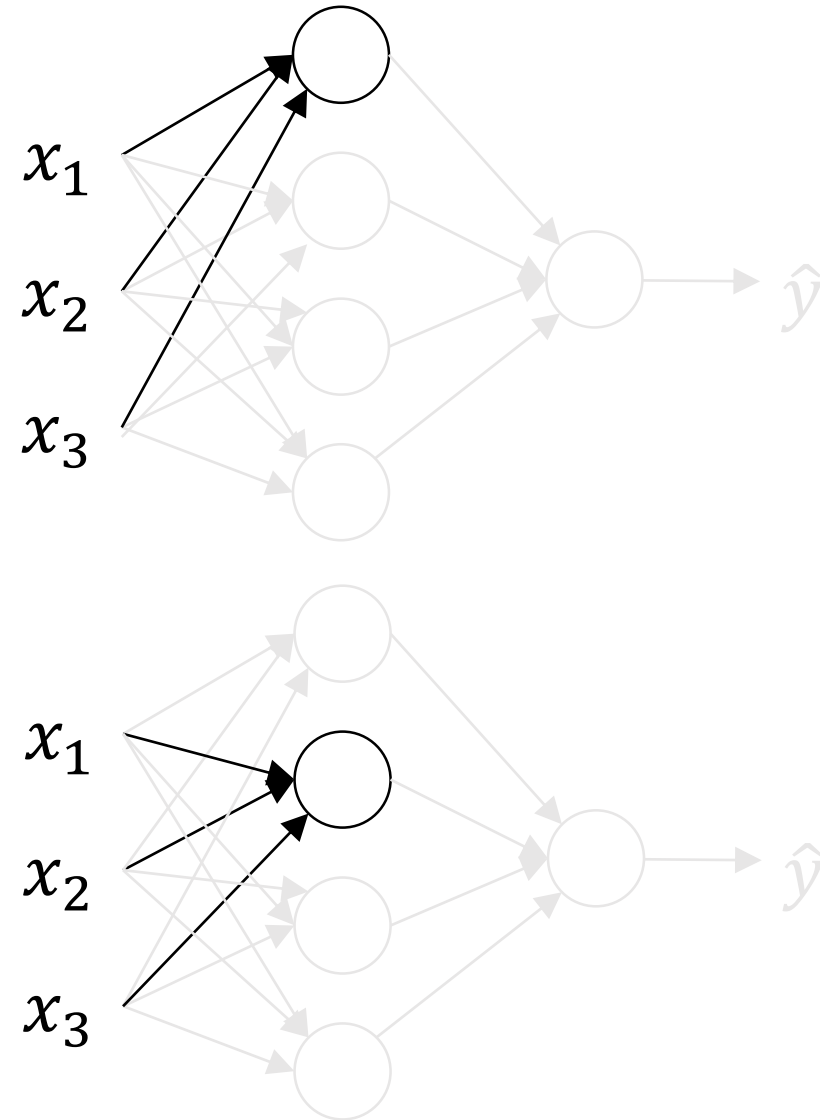


# Neural Network Representation

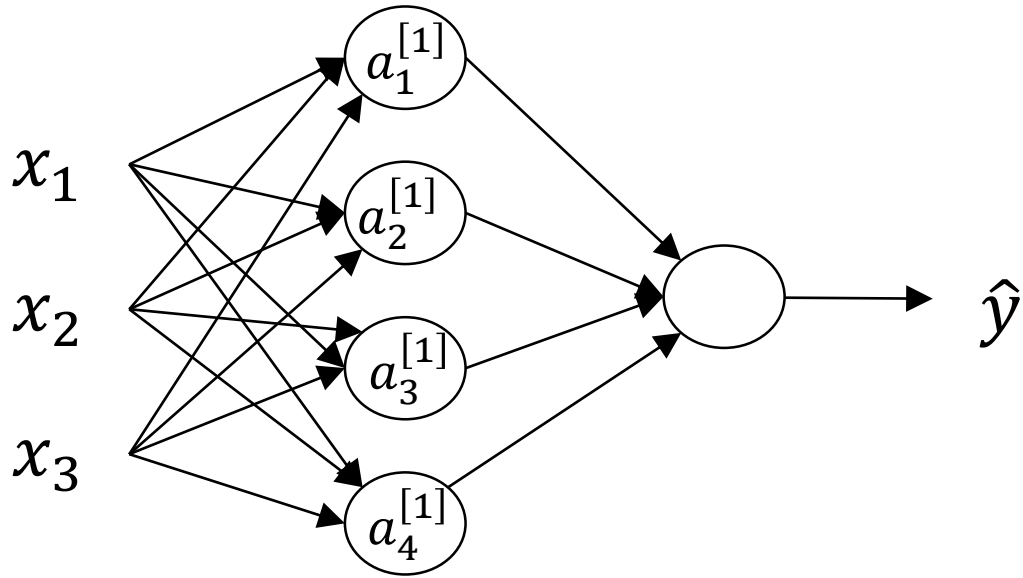


$$z = w^T x + b$$

$$a = \sigma(z)$$



# Neural Network Representation



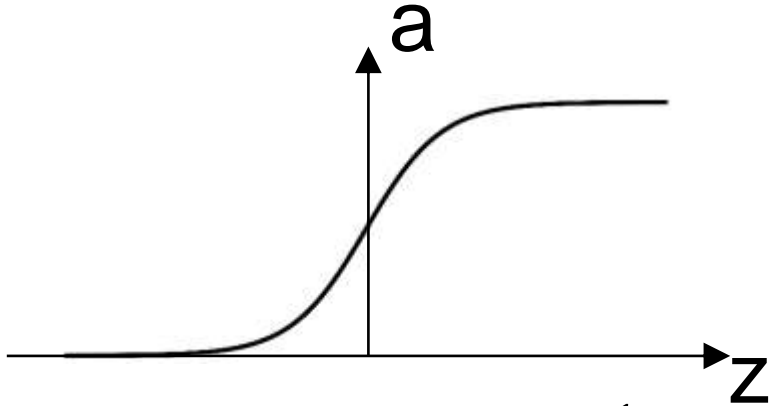
$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

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

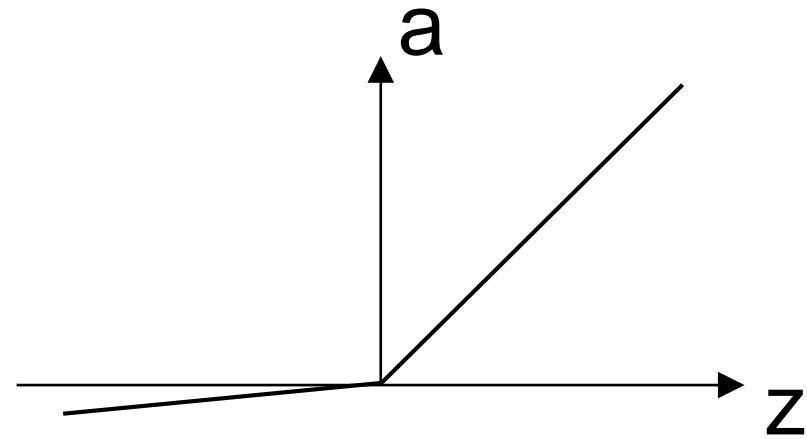
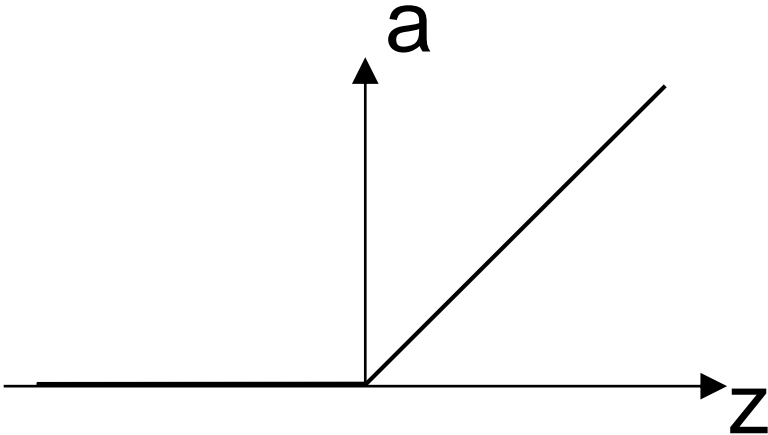
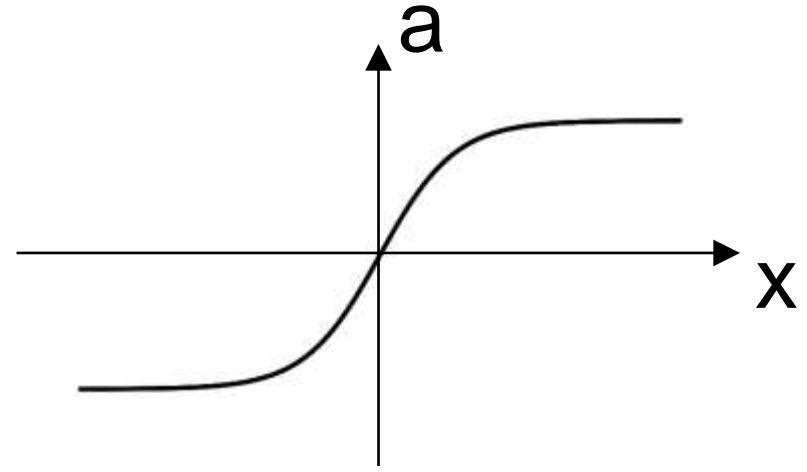
$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

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

# Activation functions



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





deeplearning.ai

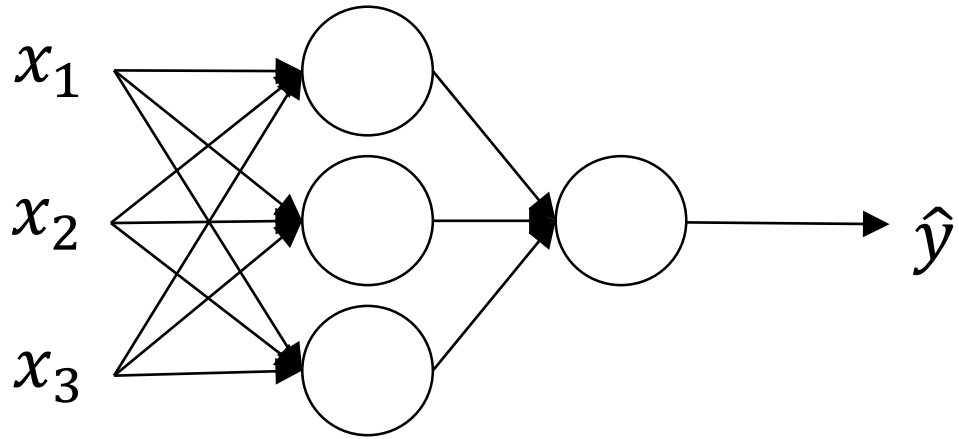
# One hidden layer Neural Network

---

Why do you  
need non-linear  
activation functions?



# Activation function



Given  $x$ :

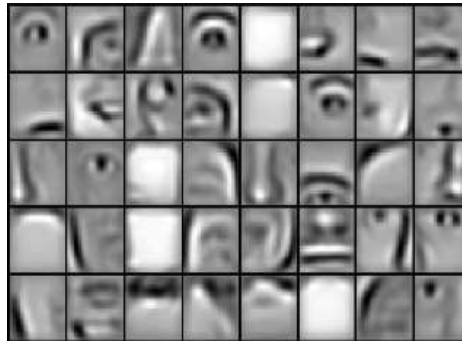
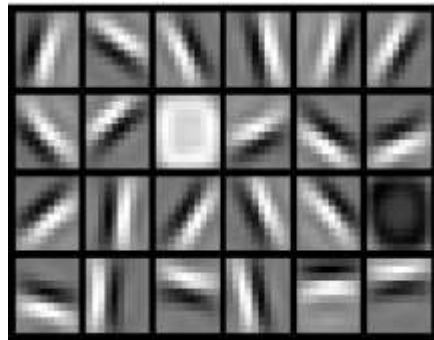
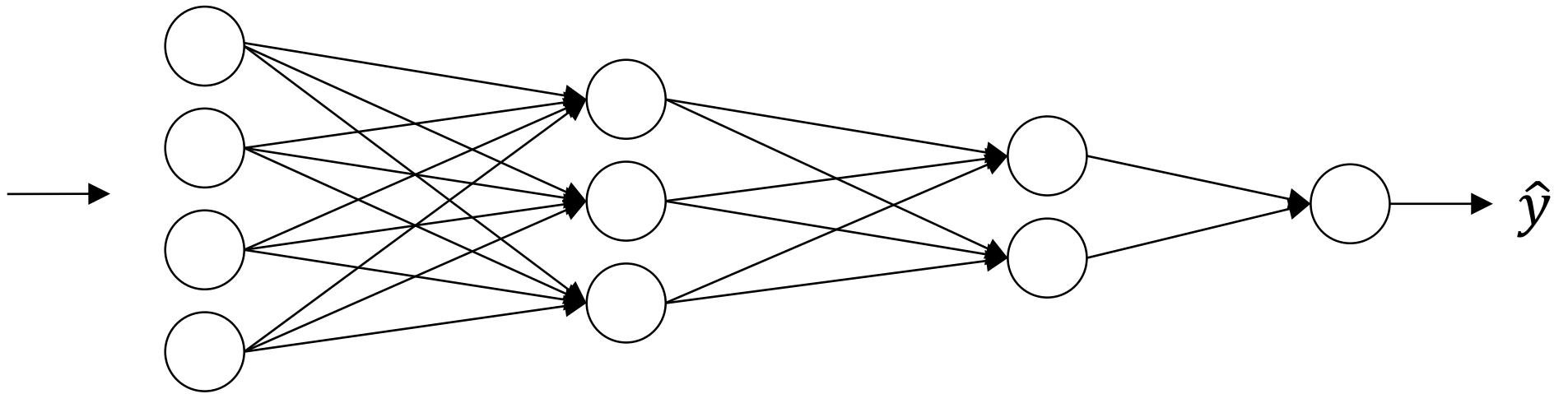
$$z^{[1]} = W^{[1]}x + b^{[1]}$$

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

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

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

# Intuition about deep representation





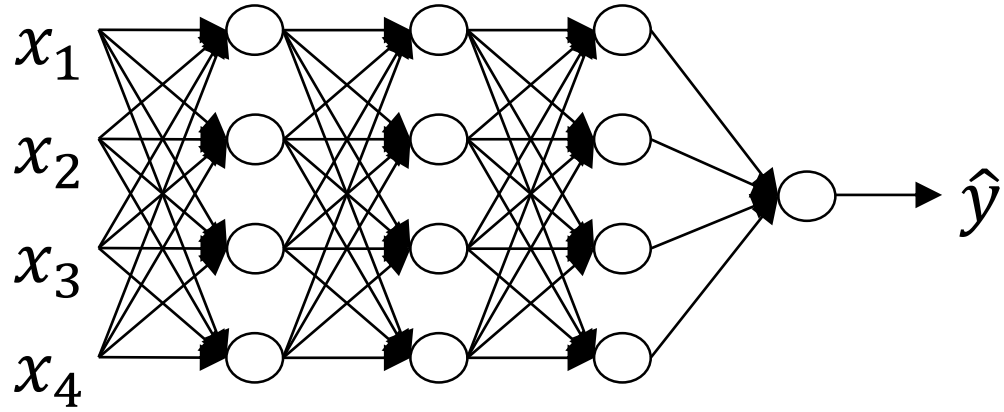
deeplearning.ai

# Deep Neural Networks

---

Building blocks of  
deep neural networks

# Forward and backward functions



# Forward and backward functions

# Forward propagation for layer $l$

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

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

# Backward propagation for layer $l$

Input  $da^{[l]}$

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

# Summary



# What are hyperparameters?

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



deeplearning.ai

# Setting up your ML application

---

## Train/dev/test sets

Train/dev/test sets

# Mismatched train/test distribution

Training set:

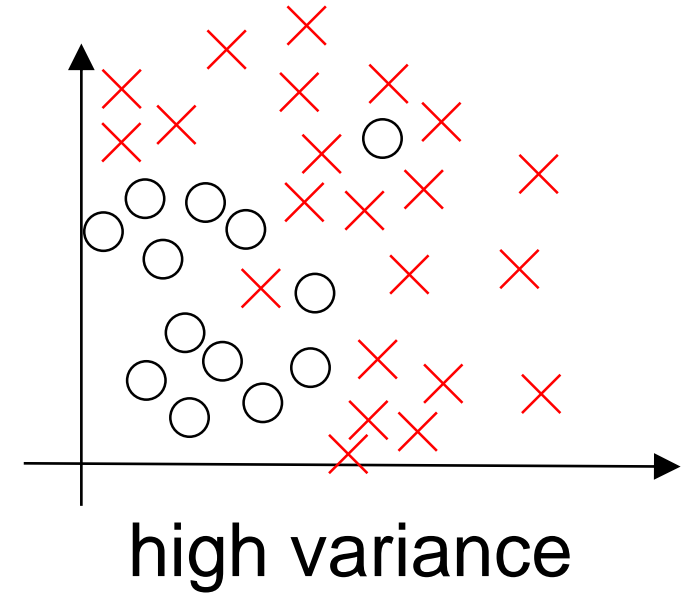
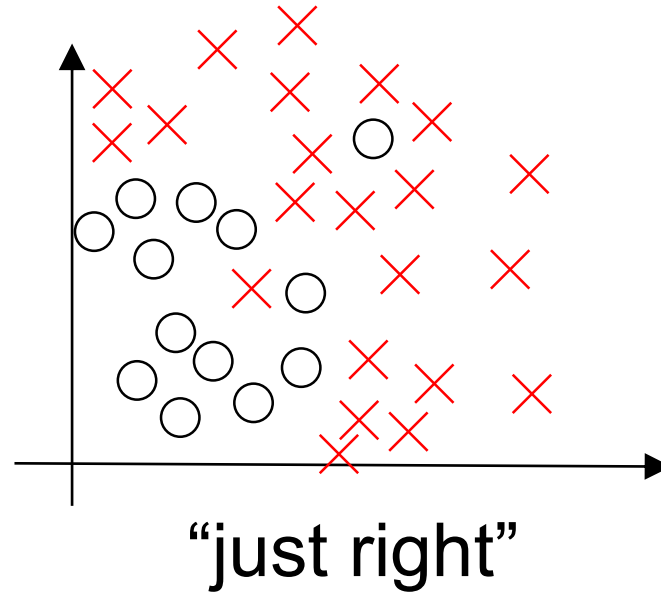
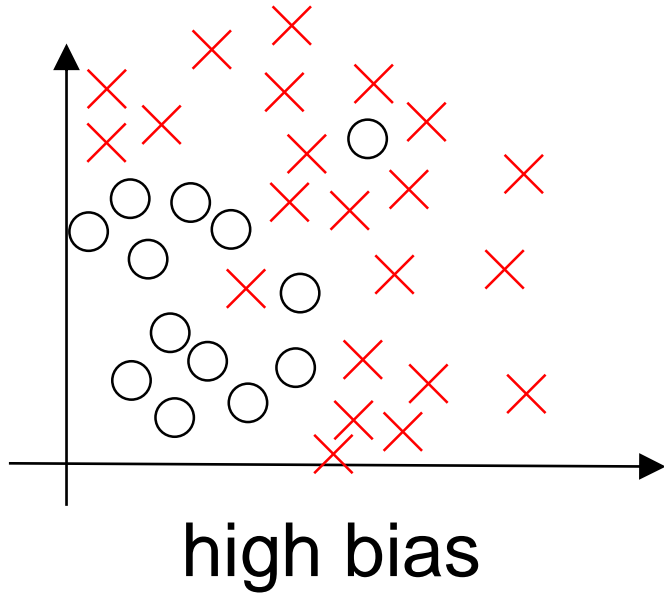
Cat pictures from  
webpages

Dev/test sets:

Cat pictures from  
users using your app

Not having a test set might be okay. (Only dev set.)

# Bias and Variance



# Bias and Variance

## Cat classification



Train set error:

Dev set error:

# Regularization techniques

- L1, L2
- Dropout
- Data Augmentation
- Early Stopping



deeplearning.ai

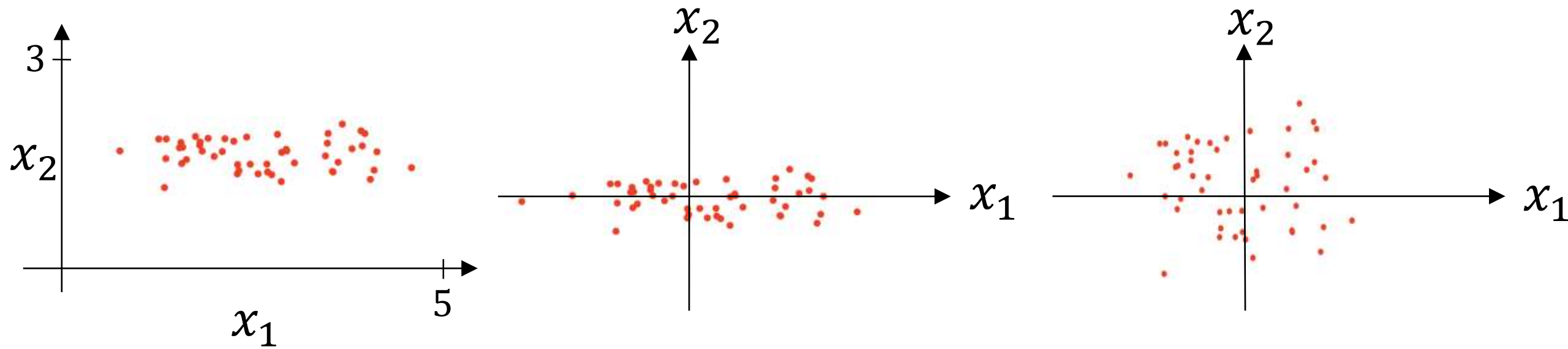
Setting up your  
optimization problem

---

Normalizing inputs



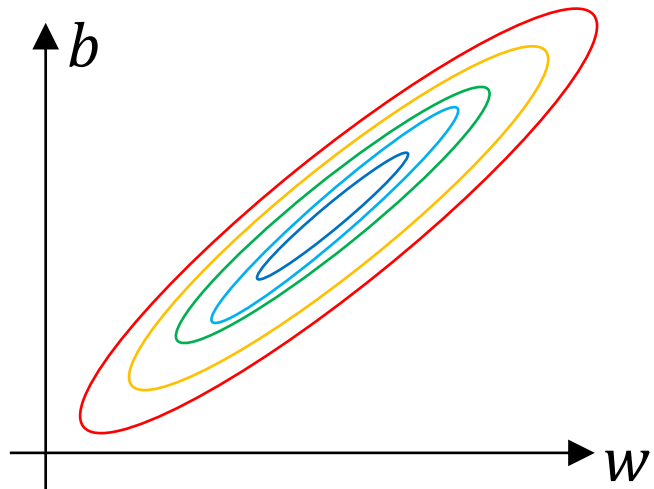
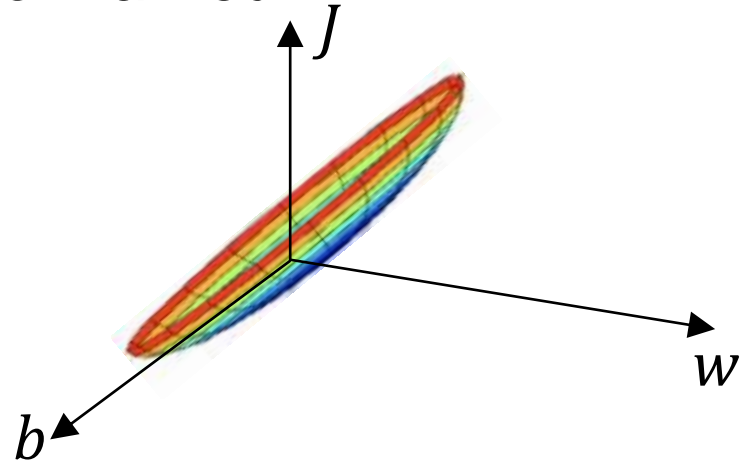
# Normalizing training sets



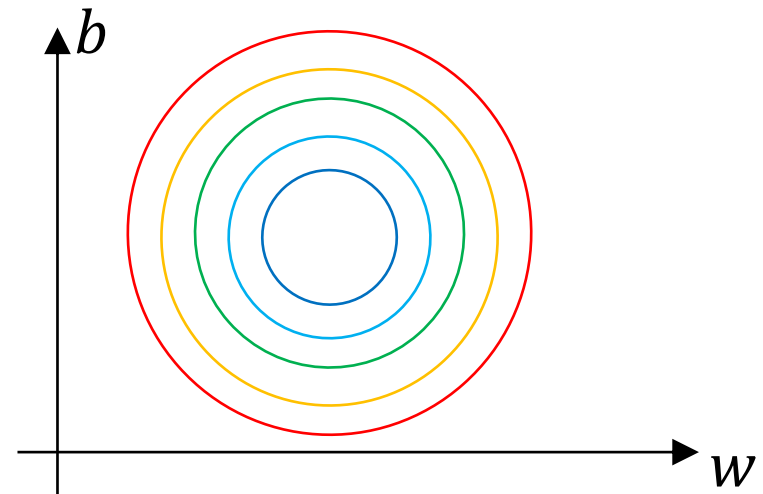
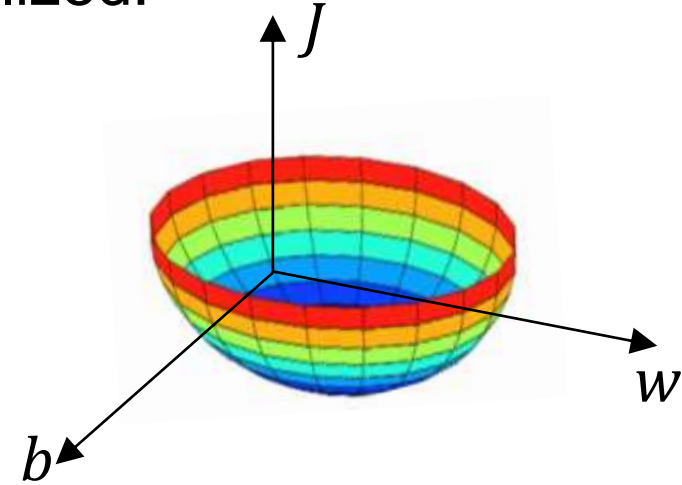
# Why normalize inputs?

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

Unnormalized:



Normalized:



# Optimization algorithms

- Gradient descent
- Mini-batch gradient descent
- Stochastic gradient descent
- RMSPROP
- ADAGRAD
- ADAM