



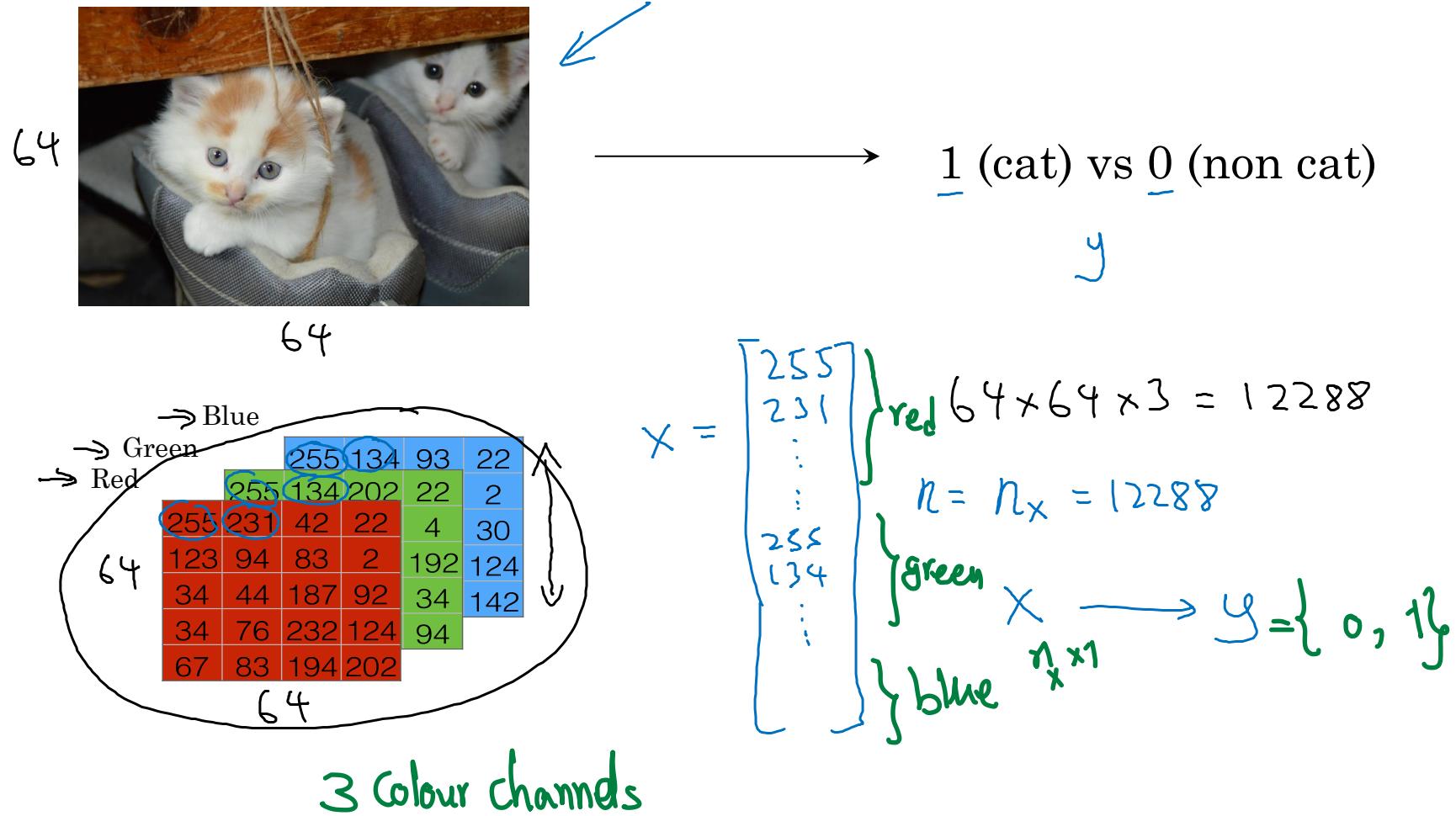
deeplearning.ai

# Basics of Neural Network Programming

---

## Binary Classification

# Binary Classification



# Notation

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

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

$$M = M_{\text{train}}$$

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

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

dim of each example

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

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

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

$$X \in \mathbb{R}^{n_x \times m}$$

$$X.\text{shape} = (n_x, m)$$

Andrew Ng

## Binary Classification

In a binary classification problem, the result is a discrete value output.

For example

- account hacked (1) or compromised (0)
- a tumor malign (1) or benign (0)

Example: Cat vs Non-Cat

The goal is to train a classifier that the input is an image represented by a feature vector,  $x$ , and predicts whether the corresponding label  $y$  is 1 or 0. In this case, whether this is a cat image (1) or a non-cat image (0).



An image is stored in the computer in three separate matrices corresponding to the Red, Green, and Blue color channels of the image. The three matrices have the same size as the image, for example, the resolution of the cat image is 64 pixels X 64 pixels, the three matrices (RGB) are 64 X 64 each.

The value in a cell represents the pixel intensity which will be used to create a feature vector of n-dimension. In pattern recognition and machine learning, a feature vector represents an object, in this case, a cat or no cat.

To create a feature vector,  $x$ , the pixel intensity values will be “unroll” or “reshape” for each color. The dimension of the input feature vector  $x$  is  $n_x = 64 \times 64 \times 3 = 12,288$ .

$n_x = \text{dim of each example}$

$m = \text{number of examples}$

$$x = \begin{bmatrix} 255 \\ 231 \\ 42 \\ \vdots \\ 255 \\ 134 \\ 202 \\ \vdots \\ 255 \\ 134 \\ 93 \\ \vdots \end{bmatrix} \longrightarrow \begin{array}{l} \text{Intensity} \\ \text{red} \\ \text{green} \\ \text{blue} \end{array}$$

$64 \times 64 \times 3 = \# \text{rows} = n_x$

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

$$y \in \mathbb{R}, \quad \gamma \in \mathbb{R}^{1 \times m} \Rightarrow \gamma.\text{shape} = (1, m)$$



deeplearning.ai

# Basics of Neural Network Programming

---

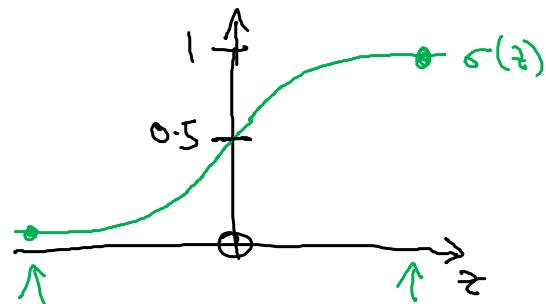
## Logistic Regression

# Logistic Regression

Given  $x$ , want  $\hat{y} = \frac{P(y=1|x)}{0 \leq \hat{y} \leq 1}$   
 $x \in \mathbb{R}^{n_x}$

Parameters :  $\underline{\omega} \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$ .

Output  $\hat{y} = \sigma(\underbrace{\omega^T x + b}_z)$



$$x_0 = 1, \quad x \in \mathbb{R}^{n_x+1}$$

$$\hat{y} = \sigma(\underline{\omega}^T x)$$

$$\underline{\omega} = \begin{bmatrix} \omega_0 \\ \omega_1 \\ \omega_2 \\ \vdots \\ \omega_{n_x} \end{bmatrix} \quad \left. \begin{array}{l} \{b \leftarrow \omega_0\} \\ \{w \leftarrow \omega_1, \dots, \omega_{n_x}\} \end{array} \right.$$

we don't use this notation.

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

If  $z$  large  $\sigma(z) \approx \frac{1}{1+0} = 1$

If  $z$  large negative number

$$\sigma(z) = \frac{1}{1 + e^{-z}} \approx \frac{1}{1 + \text{Bignum}} \approx 0$$

## Logistic Regression

Logistic regression is a learning algorithm used in a supervised learning problem when the output  $y$  are all either zero or one. The goal of logistic regression is to minimize the error between its predictions and training data.

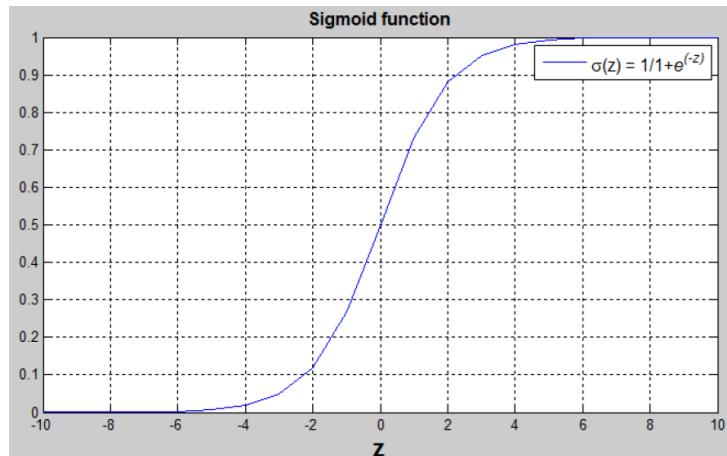
Example: Cat vs No - cat

Given an image represented by a feature vector  $x$ , the algorithm will evaluate the probability of a cat being in that image.

$$\text{Given } x, \hat{y} = P(y=1|x), \text{ where } 0 \leq \hat{y} \leq 1$$

The parameters used in Logistic regression are:

- The input features vector:  $x \in \mathbb{R}^{n_x}$ , where  $n_x$  is the number of features
- The training label:  $y \in \{0,1\}$
- The weights:  $w \in \mathbb{R}^{n_x}$ , where  $n_x$  is the number of features
- The threshold:  $b \in \mathbb{R}$
- The output:  $\hat{y} = \sigma(w^T x + b)$
- Sigmoid function:  $s = \sigma(w^T x + b) = \sigma(z) = \frac{1}{1+e^{-z}}$



$(w^T x + b)$  is a linear function ( $ax + b$ ), but since we are looking for a probability constraint between  $[0,1]$ , the sigmoid function is used. The function is bounded between  $[0,1]$  as shown in the graph above.

Some observations from the graph:

- If  $z$  is a large positive number, then  $\sigma(z) = 1$
- If  $z$  is small or large negative number, then  $\sigma(z) = 0$
- If  $z = 0$ , then  $\sigma(z) = 0.5$

$$\hat{y} = P(y=1|x) = \sigma(\underbrace{w^T x + b}_z) \Rightarrow \hat{y} = \begin{cases} 1 & z \rightarrow +\infty \\ 0 & z \rightarrow -\infty \end{cases}$$

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



deeplearning.ai

# Basics of Neural Network Programming

---

## Logistic Regression cost function

# Logistic Regression cost function

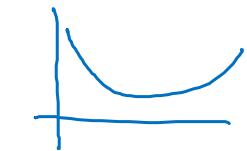
$$\rightarrow \hat{y}^{(i)} = \sigma(\underline{w^T x^{(i)}} + b), \text{ where } \sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}} \quad z^{(i)} = w^T x^{(i)} + b$$

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

$x^{(i)}$   
 $y^{(i)}$   
 $z^{(i)}$

$i$ -th example.

**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

**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

## Logistic Regression: Cost Function

To train the parameters  $w$  and  $b$ , we need to define a cost function.

Recap:

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

$x^{(i)}$  the i-th training example

Goal

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

Loss (error) function:

The loss function measures the discrepancy between the prediction ( $\hat{y}^{(i)}$ ) and the desired output ( $y^{(i)}$ ).

In other words, the loss function computes the error for a single training example.

$$L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{2}(\hat{y}^{(i)} - y^{(i)})^2 \quad \text{leads to non-convex problem X} \\ L(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})) \quad \text{and convex} \quad \checkmark \quad (\text{for Logistic Reg})$$

- If  $y^{(i)} = 1$ :  $L(\hat{y}^{(i)}, y^{(i)}) = -\log(\hat{y}^{(i)})$  where  $\log(\hat{y}^{(i)})$  and  $\hat{y}^{(i)}$  should be close to 1
- If  $y^{(i)} = 0$ :  $L(\hat{y}^{(i)}, y^{(i)}) = -\log(1 - \hat{y}^{(i)})$  where  $\log(1 - \hat{y}^{(i)})$  and  $\hat{y}^{(i)}$  should be close to 0

Cost function

The cost function is the average of the loss function of the entire training set. We are going to find the parameters  $w$  and  $b$  that minimize the overall cost function.

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

*cost*      *loss*

Cost function = the average of loss functions over the training set.

$$\begin{cases} \text{If } y=1 & P(y|x) = \hat{y} \\ \text{If } y=0 & P(y|x) = 1-\hat{y} \end{cases} \Rightarrow P(y|x) = \hat{y}^{y} (1-\hat{y})^{1-y} \Rightarrow \mathcal{L}(\hat{y}, y) = -\log P(y|x)$$

MLE  $\prod_{i=1}^m P(y^{(i)}|x^{(i)}) = P(\text{all labels})$

$$\Rightarrow \log P(\text{all labels}) = \sum_{i=1}^m \underbrace{\log P(y^{(i)}|x^{(i)})}_{-\mathcal{L}(\hat{y}^{(i)}, y^{(i)})} \Rightarrow J \equiv \frac{1}{m} \sum \mathcal{L}$$



deeplearning.ai

# Basics of Neural Network Programming

---

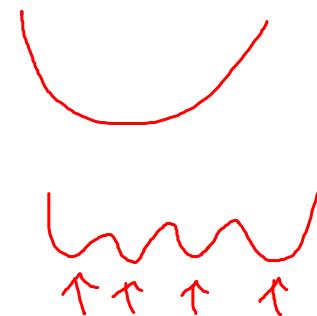
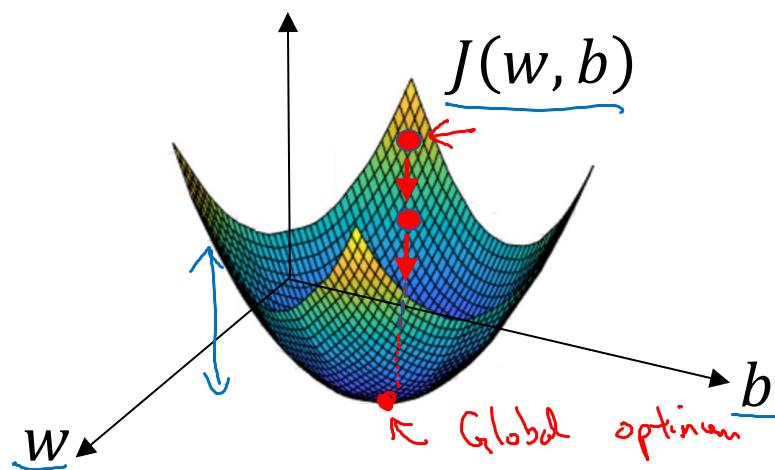
## Gradient Descent

# Gradient Descent

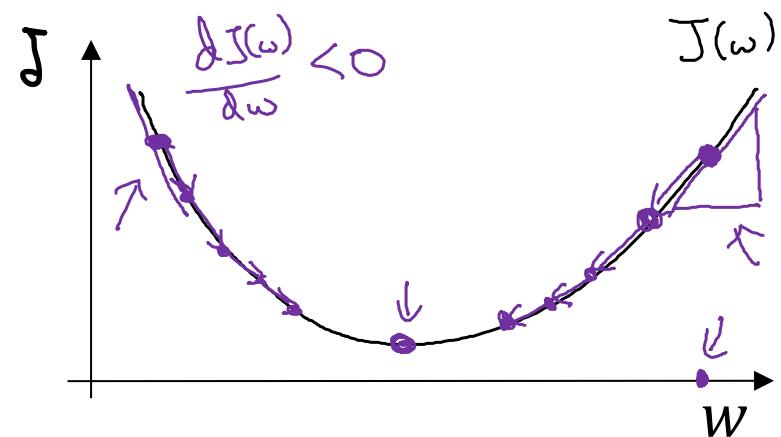
Recap:  $\hat{y} = \sigma(w^T x + b)$ ,  $\sigma(z) = \frac{1}{1+e^{-z}}$  

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

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



# Gradient Descent



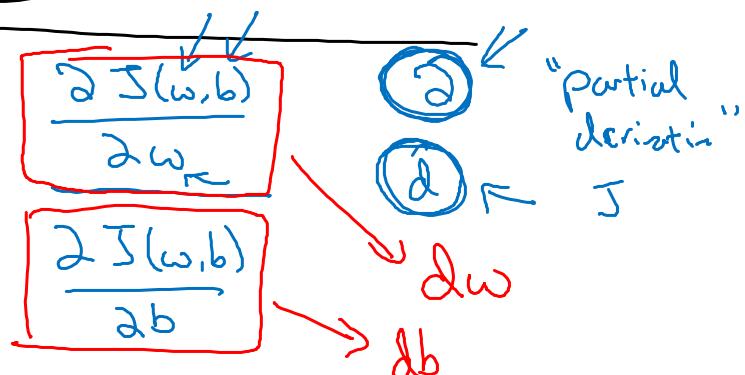
Repeat {  
 $w := w - \alpha$   
 $\uparrow \uparrow$   
 $w := w - \underline{\alpha} \underline{dw}$   
 $\frac{\partial J(w)}{\partial w} = ?$

learning rate  
"dw"

$$J(w, b)$$

$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$



$$\begin{cases} w = w - \alpha \frac{\partial J}{\partial w} \\ b = b - \alpha \frac{\partial J}{\partial b} \end{cases}$$

In Python, represent  $\frac{\partial J}{\partial w} = dw$ ,  $\frac{\partial J}{\partial b} = db$

Andrew Ng



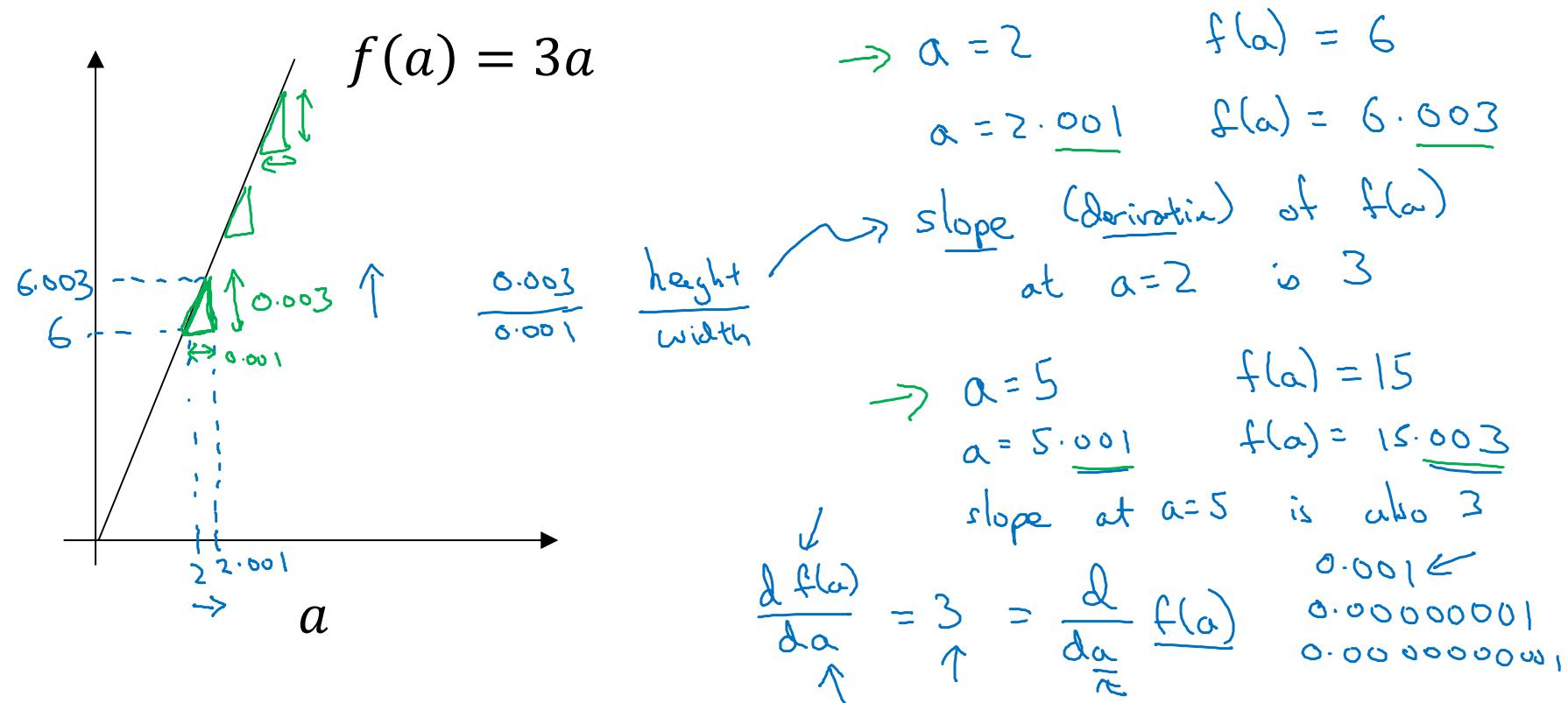
deeplearning.ai

# Basics of Neural Network Programming

---

## Derivatives

# Intuition about derivatives





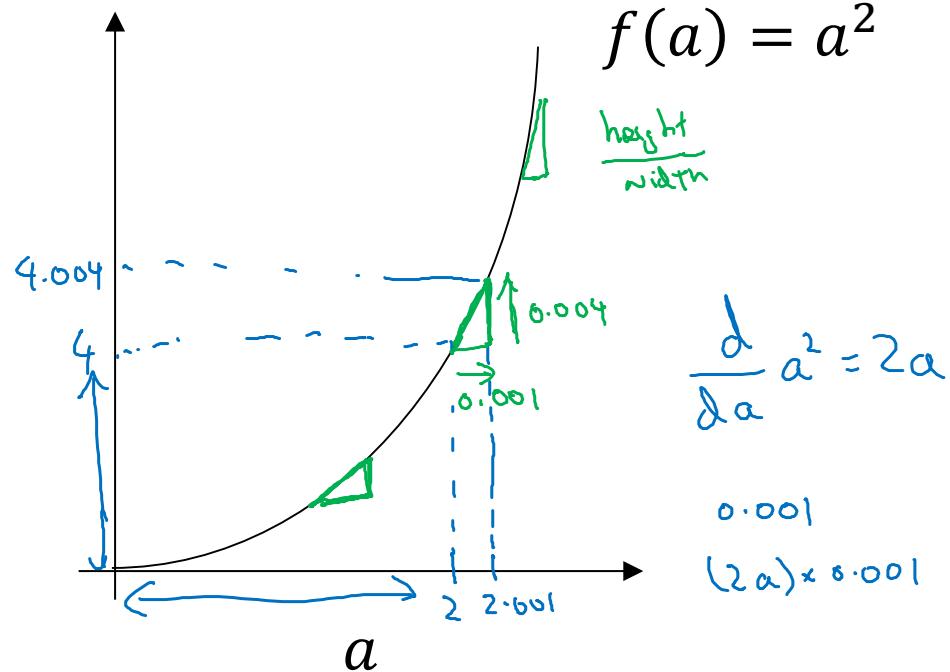
deeplearning.ai

# Basics of Neural Network Programming

---

## More derivatives examples

# Intuition about derivatives



$a=2$        $f(a)=4$   
 $a=2.001$        $f(a) \approx 4.004$   
 $\frac{(4.004 - 4)}{0.001}$

$0.001 \leftarrow$   
 $0.00000\dots 01 \leftarrow$

slope (derivative) of  $f(a)$  at  $a=2$  is 4.

$\boxed{\frac{d}{da} f(a) = 4}$  when  $\boxed{a=2}$ .

$a=5$        $f(a)=25$   
 $a=5.001$        $f(a) \approx 25.010$

$\boxed{\frac{d}{da} f(a) = 10}$  when  $\boxed{a=5}$

$\boxed{\frac{d}{da} f(a) = \frac{d}{da} a^2} = \boxed{2a}$

Andrew Ng

## More derivative examples

$$f(a) = a^2$$

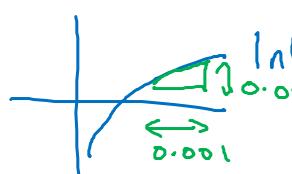
$$\frac{d}{da} f(a) = \underbrace{2a}_{4}$$

$$f(a) = a^3$$

$$\frac{d}{da} f(a) = \underbrace{3a^2}_{3+2^2 = 12}$$

$$f(a) = \frac{\log_e(a)}{\ln(a)}$$

$$\frac{d}{da} f(a) = \frac{1}{a}$$


 $\frac{d}{da} f(a) = \boxed{\frac{1}{2}}$

$$a = 2$$

$$a = 2.001$$

$$f(a) = 4$$

$$f(a) \approx 4.004$$

$$a = 2$$

$$a = \underline{2.001}$$

$$f(a) = 8$$

$$f(a) \approx \underline{8.012}$$

$$a = 2$$

$$a = \underline{2.001}$$

$$f(a) \approx 0.69315$$

$$f(a) \approx \underline{0.69365}$$

$$0.0065 \xrightarrow{0.0005}$$

Andrew Ng



deeplearning.ai

# Basics of Neural Network Programming

---

## Computation Graph

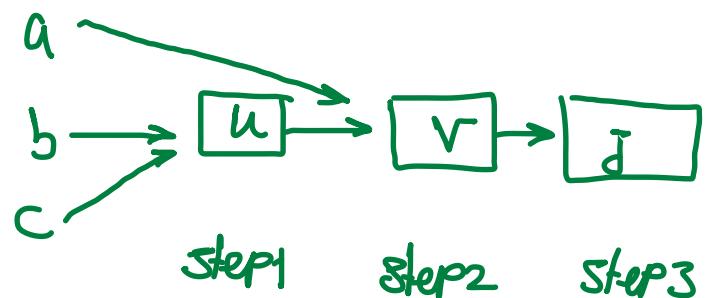
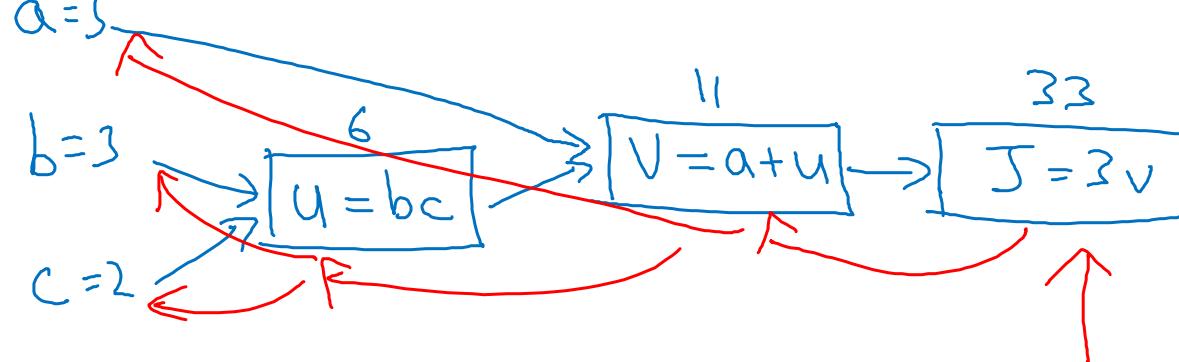
# Computation Graph

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

$\underbrace{\phantom{000}}_{J}$

3 steps

$$\begin{cases} J(a, b, c) = 3(a + bc) \\ u = bc \\ v = a + u \\ J = 3v \end{cases}$$





deeplearning.ai

# Basics of Neural Network Programming

---

## Derivatives with a Computation Graph

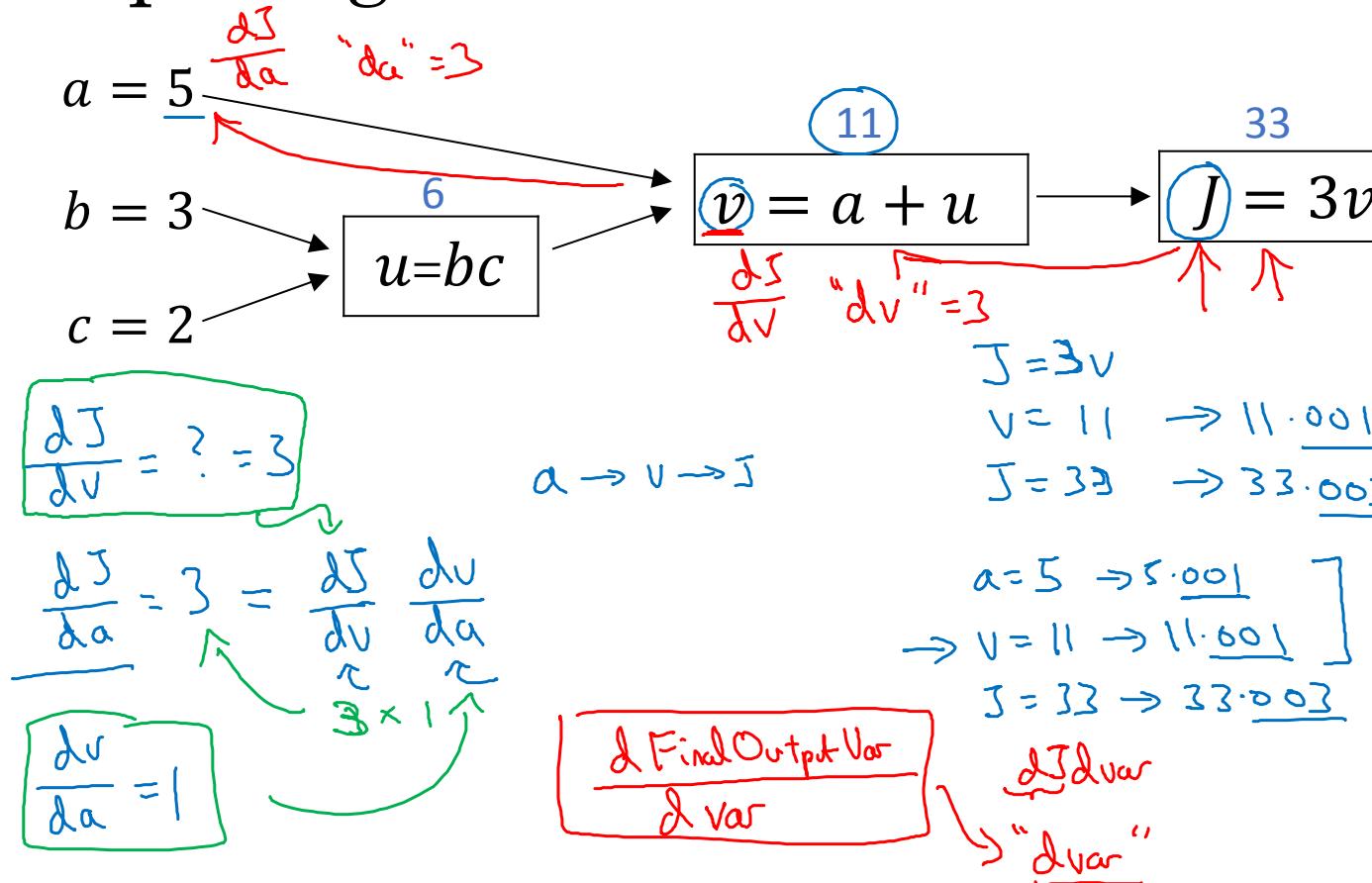
Goal:  $\frac{\partial(\text{Final output variable})}{\partial \text{Var}} = "d\text{var}" \equiv \frac{\partial J}{\partial \text{Var}} = ?$

↓

Back propagation  $d\alpha = \frac{\partial J}{\partial \alpha} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial \alpha} = dv \frac{\partial v}{\partial \alpha}$

Python:  $du = \frac{\partial J}{\partial u}, d\alpha = \frac{\partial J}{\partial \alpha}, \dots$

## Computing derivatives



$$f(a) = 3a$$

$$\frac{df(w)}{da} = \frac{df}{da} = 3$$

$$J = 3v$$

$$\frac{dJ}{dv} = 3$$

# Computing derivatives

$$\begin{aligned}
 \frac{\partial J}{\partial a} &\rightarrow a = 5 \quad \frac{\partial a}{\partial a} = 1 \\
 \frac{\partial J}{\partial b} &\rightarrow b = 3 \quad \frac{\partial b}{\partial b} = 1 \\
 \frac{\partial J}{\partial c} &\rightarrow c = 2 \quad \frac{\partial c}{\partial c} = 1
 \end{aligned}$$

$u = bc$

$$\begin{array}{c}
 \text{11} \\
 v = a + u \\
 \frac{\partial v}{\partial v} = 1
 \end{array}
 \quad
 \begin{array}{c}
 33 \\
 J = 3v \\
 \frac{\partial J}{\partial J} = 1
 \end{array}$$

$$\frac{\partial J}{\partial u} = 3 = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} = 3 \cdot 1$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial u} \cdot \frac{\partial u}{\partial b} = 3 \cdot 2 = 6$$

$$\frac{\partial J}{\partial a} = \frac{\partial J}{\partial u} \cdot \frac{\partial u}{\partial a} = 3 \cdot 3 = 9$$

$u = 6 \rightarrow 6.001$   
 $v = 11 \rightarrow 11.001$   
 $J = 33 \rightarrow 33.003$

$b = 3 \rightarrow 3.001$   
 $u = b \cdot c = 6 \rightarrow 6.002$   
 $J = 33.006$

$v = 11.002$   
 $J = 33.006$

Andrew Ng



deeplearning.ai

# Basics of Neural Network Programming

---

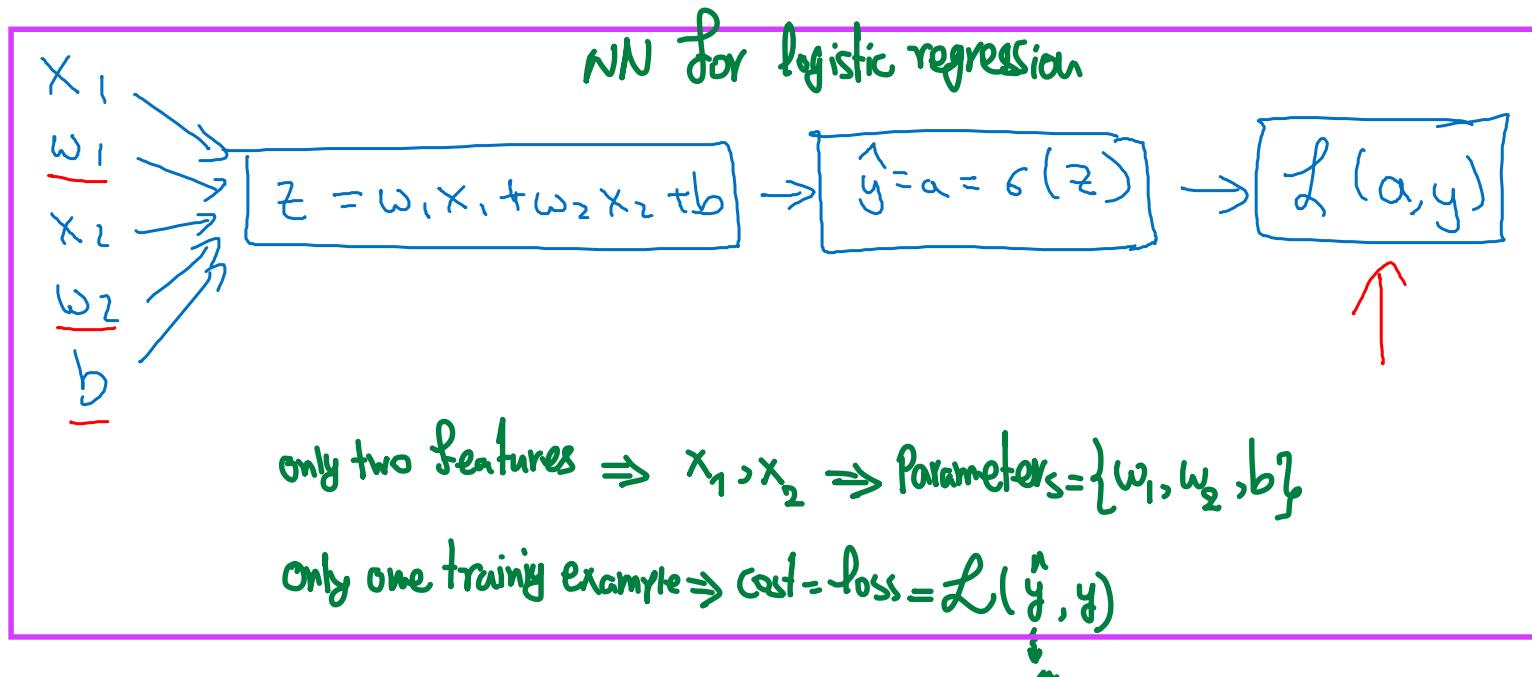
## Logistic Regression Gradient descent

# Logistic regression recap

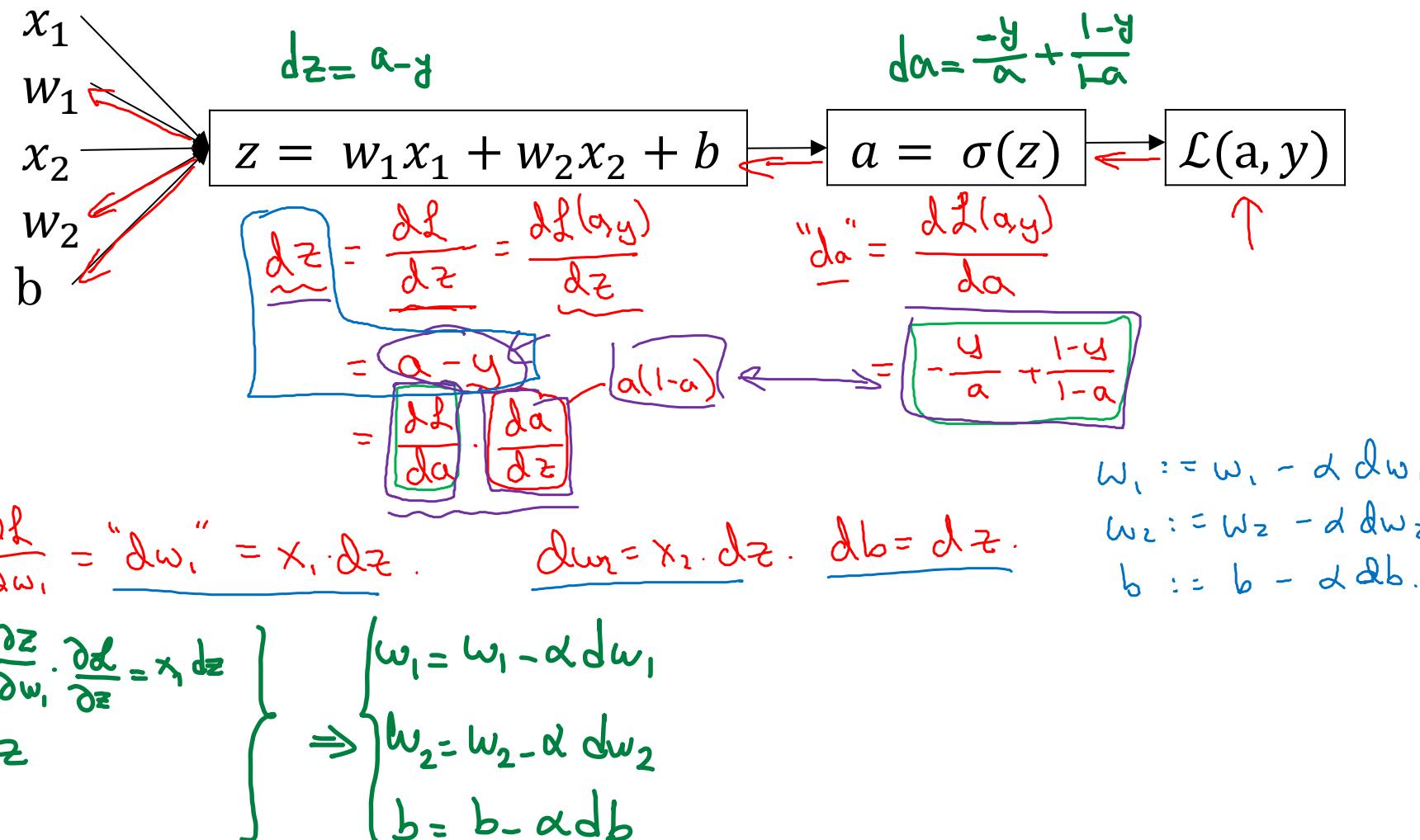
$$\rightarrow z = w^T x + b$$

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

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



# Logistic regression derivatives





deeplearning.ai

# Basics of Neural Network Programming

---

Gradient descent  
on  $m$  examples

# Logistic regression on $m$ examples

$$\begin{aligned} J(\omega, b) &= \underbrace{\frac{1}{m} \sum_{i=1}^m \ell(a^{(i)}, y^{(i)})}_{\text{for one observation}} \\ \rightarrow a^{(i)} &= \hat{y}^{(i)} = g(z^{(i)}) = g(\omega^\top x^{(i)} + b) \end{aligned}$$

$(x^{(i)}, y^{(i)})$

$\underline{dw_1}^{(i)}, \underline{dw_2}^{(i)}, \underline{db}^{(i)}$

$$\underline{\frac{\partial}{\partial w_1} J(\omega, b)} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_1} \ell(a^{(i)}, y^{(i)})}_{\underline{dw_1}^{(i)}} \leftarrow (x^{(i)}, y^{(i)}) \text{ for one observation}$$

$$\Rightarrow \frac{\partial}{\partial w_1} J = \frac{1}{m} \sum_{i=1}^m \underline{dw_1}^{(i)} \quad \text{where} \quad \underline{dw_1}^{(i)} = \frac{\partial \ell}{\partial w_1^{(i)}}$$

# Logistic regression on $m$ examples

loop<sub>1</sub>  
| samples)

$$J = 0; \underline{\Delta w_1} = 0; \underline{\Delta w_2} = 0; \underline{\Delta b} = 0$$

For  $i = 1$  to  $m$

$$z^{(i)} = \omega^\top x^{(i)} + b$$

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

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

$$\underline{\Delta z^{(i)}} = \alpha^{(i)} - y^{(i)}$$

$$\begin{aligned} \Delta w_1 &+= x_1^{(i)} \Delta z^{(i)} \\ \Delta w_2 &+= x_2^{(i)} \Delta z^{(i)} \end{aligned}$$

$\uparrow$        $\uparrow$   
 $n=2$   
# features

$$\begin{aligned} \Delta b &+= \Delta z^{(i)} \\ \Delta w_3 & \\ \vdots & \\ \Delta w_n & \end{aligned}$$

$$J /= m \leftarrow$$

$$\Delta w_1 /= m; \Delta w_2 /= m; \Delta b /= m. \leftarrow$$

$\uparrow$        $\uparrow$        $\uparrow$   
average over  $m$  samples

$$\Delta w_1 = \frac{\partial J}{\partial w_1}$$

$$w_1 := w_1 - \alpha \underline{\Delta w_1}$$

$$w_2 := w_2 - \alpha \underline{\Delta w_2}$$

$$b := b - \alpha \underline{\Delta b}.$$

Vectorization  $\rightarrow$

To prevent the for loops



deeplearning.ai

# Basics of Neural Network Programming

---

## Vectorization

# What is vectorization?

$$z = \underbrace{w^T x}_{} + b$$

Non-vectorized:

```
z = 0  
for i in range(n_x):  
    z += w[i] * x[i]
```

$z += b$

$$w = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$
$$w \in \mathbb{R}^{n_x}$$
$$x \in \mathbb{R}^{n_x}$$

Vectorized

$$z = \underbrace{\text{np.dot}(w, x)}_{w^T x} + b$$

$\rightarrow$  GPU } SIMD - single instruction  
 $\rightarrow$  CPU } multiple data.

Vectorization demo

jupyter Vectorization demo Last Checkpoint: 3 minutes ago (unsaved changes)

Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

```
a = np.random.rand(1000000)
b = np.random.rand(1000000)

tic = time.time()
c = np.dot(a,b)
toc = time.time()

print(c)
print("Vectorized version:" + str(1000*(toc-tic)) + "ms")

c = 0
tic = time.time()
for i in range(1000000):
    c += a[i]*b[i]
toc = time.time()

print(c)
print("For loop:" + str(1000*(toc-tic)) + "ms")
```

250286.989866  
Vectorized version:1.5027523040771484ms  
250286.989866  
For loop:474.29513931274414ms

vectorized → 1,5 ms  
non-vectorized → 474,2 ms



deeplearning.ai

# Basics of Neural Network Programming

---

## More vectorization examples

# Neural network programming guideline

Whenever possible, avoid explicit for-loops.

non-vectorized

$$u = Av$$

$$u_i = \sum_i \sum_j A_{ij} v_j$$

$$u = np.zeros((n,))$$

for i ...  $\leftarrow$

for j ...  $\leftarrow$

$$u[i] += A[:, i] * v[i]$$

vectorized

$$u = np.dot(A, v)$$

Andrew Ng

# Vectors and matrix valued functions

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

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

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

non-vectorized

```
import numpy as np  
u = np.exp(v) ←  
→  
np.log(v)  
np.abs(v)  
np.maximum(v, 0)  
v**2           /v
```

vectorized operations in numpy

exp, log, abs, maximum, \*\* , / ,

Samples مربوط

Features مربوط

نواته خذ گند ← جلوه همچنان با هم خویش گند.

## Logistic regression derivatives

$$J = 0,$$

$$\boxed{dw_1 = 0, dw_2 = 0}, \quad db = 0 \quad \Rightarrow \quad dw = np.zeros((n_x, 1))$$

for i = 1 to n:

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

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

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

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

for i=1...n  
dw<sub>j</sub> = ...  
x<sup>(i)</sup> = ...

$$\boxed{dw_1 += x_1^{(i)} dz^{(i)}}, \quad n_x = 2$$

$$\boxed{dw_2 += x_2^{(i)} dz^{(i)}} \quad \Rightarrow \quad dw += x^{(i)} dz^{(i)}$$

$$J = J/m, \quad \boxed{dw_1 = dw_1/m, \quad dw_2 = dw_2/m}, \quad db = db/m$$

$$\Rightarrow dw /= m$$



deeplearning.ai

# Basics of Neural Network Programming

---

## Vectorizing Logistic Regression

# Vectorizing Logistic Regression

$$\begin{aligned} \rightarrow z^{(1)} &= w^T x^{(1)} + b \\ \rightarrow a^{(1)} &= \sigma(z^{(1)}) \end{aligned}$$

$$\begin{aligned} \rightarrow z^{(2)} &= w^T x^{(2)} + b \\ \rightarrow a^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

$$\begin{aligned} \rightarrow z^{(3)} &= w^T x^{(3)} + b \\ \rightarrow a^{(3)} &= \sigma(z^{(3)}) \end{aligned}$$

$$\underline{\underline{X}} = \left[ \begin{array}{c|c|c|c} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \hline | & | & & | \\ \hline \end{array} \right]$$

$$\frac{(n_x, m)}{\mathbb{R}^{n_x \times m}}$$

$$\underline{\underline{w}}^\top \left[ \begin{array}{c|c|c|c} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \hline | & | & & | \\ \hline \end{array} \right]$$

$$\underline{\underline{z}} = \left[ \begin{array}{c|c|c|c} z^{(1)} & z^{(2)} & \dots & z^{(m)} \\ \hline | & | & & | \\ \hline \end{array} \right] = \underline{\underline{w}}^\top \underline{\underline{X}} + \left[ \begin{array}{c|c|c|c} b & b & \dots & b \\ \hline | & | & & | \\ \hline \end{array} \right] = \left[ \begin{array}{c|c|c|c} w^T x^{(1)} + b & w^T x^{(2)} + b & \dots & w^T x^{(m)} + b \\ \hline | & | & & | \\ \hline \end{array} \right]$$

$$\rightarrow \underline{\underline{z}} = \text{np.dot}(\underline{\underline{w}}^\top, \underline{\underline{X}}) + \underline{\underline{b}}$$

"Broadcasting"

$$\underline{\underline{A}} = \left[ \begin{array}{c|c|c|c} a^{(1)} & a^{(2)} & \dots & a^{(m)} \\ \hline | & | & & | \\ \hline \end{array} \right] = \underline{\underline{\sigma}}(\underline{\underline{z}})$$

# Vectorizing Logistic Regression

$$\begin{aligned} \rightarrow z^{(1)} &= w^T x^{(1)} + b \\ \rightarrow a^{(1)} &= \sigma(z^{(1)}) \end{aligned}$$

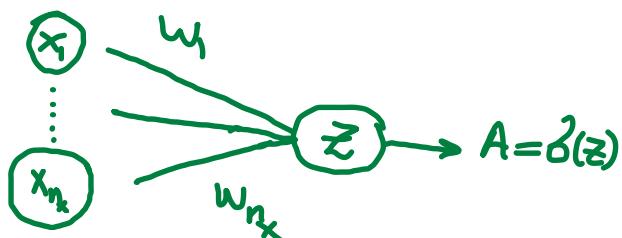
$$\begin{aligned} \rightarrow z^{(2)} &= w^T x^{(2)} + b \\ \rightarrow a^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

$$\begin{aligned} \rightarrow z^{(3)} &= w^T x^{(3)} + b \\ \rightarrow a^{(3)} &= \sigma(z^{(3)}) \end{aligned}$$

$$X = \left[ \begin{array}{c:c:c:c} \vdots & \vdots & \vdots & \vdots \\ x^{(1)} & \dots & x^{(m)} & \\ \vdots & & \vdots & \\ x^{(n)} & & & \end{array} \right]_{n \times m} \Rightarrow Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}]_{1 \times m}, \ b = [b \ b \ \dots \ b]_{1 \times m}, \ A = [a^{(1)} \ \dots \ a^{(m)}]$$

$$Z = w^T X + b \Rightarrow Z = np.\dot{dot}(w.T, X) + b$$

↓  
transpose      ↗ = broadcasting operator



$$A = \sigma(Z)$$



deeplearning.ai

# Basics of Neural Network Programming

---

## Vectorizing Logistic Regression's Gradient Computation

# Vectorizing Logistic Regression

$$dz = [dz^{(1)} \dots dz^{(m)}], \quad A = [a^{(1)} \dots a^{(m)}], \quad Y = [y^{(1)} \dots y^{(m)}]$$

$$\Rightarrow dz = A - Y$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} = \frac{1}{m} np \cdot \text{sum}(dz)$$

$$dw = \frac{1}{m} X \times dz^T = \frac{1}{m} np \cdot \text{dot}(X, dz.T)$$



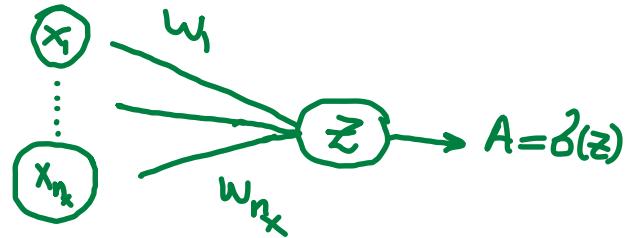
# Vectorizing Logistic Regression

$$\begin{aligned}
 dz^{(1)} &= a^{(1)} - y^{(1)} & dz^{(2)} &= a^{(2)} - y^{(2)} & \dots \\
 dz &= [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}]_{{1 \times m}} \quad \leftarrow \\
 A &= [a^{(1)} \ \dots \ a^{(m)}] \quad Y = [y^{(1)} \ \dots \ y^{(m)}] \\
 \rightarrow dz &= A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots]
 \end{aligned}$$

$\frac{\partial w_1}{\partial w_1}, \frac{\partial w_2}{\partial w_2}, \dots$   
 $\frac{\partial b}{\partial w_1}, \frac{\partial b}{\partial w_2}, \dots$   
 $\frac{\partial b}{\partial b} = m$   
 $\frac{\partial w}{\partial w} = m$   
 $\frac{\partial w}{\partial w} = \frac{dw}{m}$

$$\begin{aligned}
 db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\
 &= \frac{1}{m} \text{np.sum}(dz) \\
 dw &= \frac{1}{m} X dz^T \\
 &= \frac{1}{m} \begin{bmatrix} x^{(1)} & \dots & x^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix} \\
 &= \frac{1}{m} \left[ \underline{x^{(1)} dz^{(1)}} + \dots + \underline{x^{(m)} dz^{(m)}} \right]_{n \times 1}
 \end{aligned}$$

Andrew Ng



# Implementing Logistic Regression

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

for i = 1 to m:

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

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

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

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

$$\left\{ \begin{array}{l} dw_1 = x_1^{(i)} dz^{(i)} \\ dw_2 = x_2^{(i)} dz^{(i)} \end{array} \right\} \quad \partial w + = X^{(i)} * d\bar{z}^{(i)}$$

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

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

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

$$db = db/m$$

```

for iter in range(1000):
    Z = w^T X + b
    A = g(Z)
    dZ = A - Y
    dw = 1/m * X^T dZ^T
    db = 1/m * np.sum(dZ)

    w := w - alpha dw
    b := b - alpha db
}

```

} for Propagation

} back propagation

# Full Algorithm Gradient descent for logistic regression



deeplearning.ai

# Basics of Neural Network Programming

---

## Broadcasting in Python

۳- `Reshape` فرمان اجازه ماتریس را به دلخواه ما تغییر کند: `(A).reshape(4,4)` ابعاد  $4 \times 4$  حاکم است.

# Broadcasting example

## Calories from Carbs, Proteins, Fats in 100g of different foods:

	Apples	Beef	Eggs	Potatoes	
Carb	56.0	0.0	4.4	68.0	
Protein	1.2	104.0	52.0	8.0	= A
Fat	1.8	135.0	99.0	0.9	(3,4)
	59 cal	$\frac{56}{59} \approx 94.9\%$			

Calculate % of colors from Cub, Pwter, Fart. Can you do this without explicit for-loop?

```
cal = A.sum(axis = 0)  
percentage = 100*A/(cal.reshape(1,4))  
                                ↑(3,4) / (1,4)
```

→ broadcasting

# Broadcasting example

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} \xrightarrow{100} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m,n)} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}_{(1,n) \rightsquigarrow (m,n)} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix} \quad \downarrow \quad \downarrow \quad \downarrow$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m,n)} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix}_{\substack{(m,1) \\ (m,n)}} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix} \quad \leftarrow \quad \leftarrow$$

# General Principle of BroadCasting

$$\begin{array}{ccc} (m, n) & \xrightarrow{\quad \cancel{*} \quad} & \begin{array}{c} (1, n) \\ \hline (m, 1) \end{array} \xrightarrow{\quad \cancel{*} \quad} (m, n) \end{array}$$

$$\begin{array}{ccc} (m, l) & + & \mathbb{R} \\ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} & + & 100 \\ \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} & + & 100 \end{array} = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix} = [101 \quad 102 \quad 103]$$

Matlab/Octave: `bssxfun`  $\longleftrightarrow$  Broadcasting in MATLAB and Octave.



deeplearning.ai

# Basics of Neural Network Programming

---

A note on python/  
numpy vectors

# Python / numpy vectors

```
import numpy as np
```

(5, )	a = np.random.randn(5)	'rank 1 array'	x
(5, 1)	a = np.random.randn((5, 1))	column vector	✓
(1, 5)	a = np.random.randn((1, 5))	row vector	✓
T/F	assert(a.shape == (5, 1))	Checking dimensions	
(5, 1)	a=a.reshape((5, 1))	Transform dimensions	