# Copyright Notice

Basics of Neural Network Programming

Binary Classification

deeplearning.ai

# Binary Classification



64

64

$\longrightarrow$   1 (cat) vs 0 (non cat)

$y$

$\rightarrow$ Blue
$\rightarrow$ Green
$\rightarrow$ Red

| 255 | 134 | 93 | 22 |
|-----|-----|-----|-----|
| 255 | 134 | 202 | 22 | 2 |
| 255 | 231 | 42 | 22 | 4 | 30 |
| 123 | 94 | 83 | 2 | 192 | 124 |
| 34 | 44 | 187 | 92 | 34 | 142 |
| 34 | 76 | 232 | 124 | 94 |
| 67 | 83 | 194 | 202 |

64

64

$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix}$$

$64 \times 64 \times 3 = 12288$

$n = n_x = 12288$

$x \longrightarrow y$

Andrew Ng

# Notation

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

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

$m = m_{train}$          $m_{test} = \#test$ examples.

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

$\underset{\longleftarrow m \longrightarrow}{}$

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

$$\begin{matrix} x^{(1)} \\ \vdots \\ x^{(m)} \end{matrix}$$ (crossed out)

$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \ldots & y^{(m)} \end{bmatrix}$

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

$Y.shape = (1, m)$
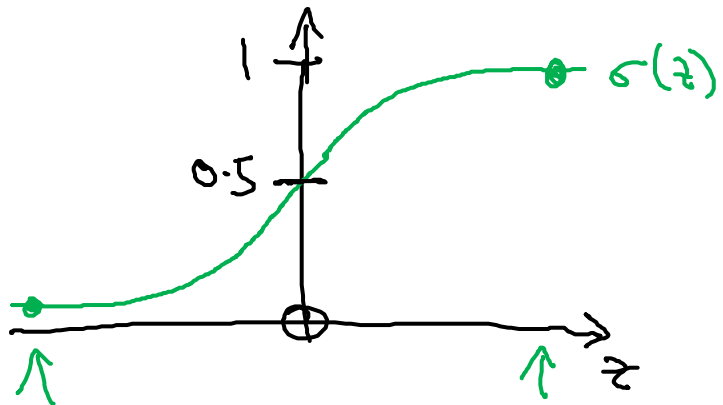
Basics of Neural
Network Programming

Logistic Regression

deeplearning.ai

# Logistic Regression

Given $x$, want $\hat{y} = P(y=1|x)$

$0 \leq \hat{y} \leq 1$

$x \in \mathbb{R}^{n_x}$

Parameters: $\boxed{w} \in \mathbb{R}^{n_x}$, $\boxed{b} \in \mathbb{R}$.

where b is the intercept

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



Logistic regression utilises the sigmoid function to fulfill probability bounds

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

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \begin{matrix} \}b \leftarrow \\ \\ \}w \leftarrow \\ \\ \end{matrix}$$

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

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

If $z$ large negative number

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

Andrew Ng

# Logistic Regression cost function

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

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

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

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

$i-th$ example.

Loss (error) function:
$\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$

$\mathcal{L}(\hat{y}, y) = -\left[ y \log \hat{y} + (1-y) \log(1-\hat{y}) \right] \leftarrow$

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

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

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

Andrew Ng

# Basics of Neural Network Programming

## Gradient Descent

deeplearning.ai

# Gradient Descent

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

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

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



$J(w, b)$

Initialisation point, and then the algo aims to converge to the global optimum minima

$b$

$\leftarrow$ Global optimum

$w$

Andrew Ng

# Gradient Descent



$\dfrac{dJ(\omega)}{d\omega} < 0$

derivative is negative, so w is increased

derivative is positive, so w is decreased

$J(\omega)$

$W$

Repeat {

$\omega := \omega - \alpha \dfrac{dJ(\omega)}{d\omega}$

learning rate

}

"dw"

$\omega := \omega - \alpha dw$

learning rate x the derivative is the update to parameters W (and b)

$\dfrac{dJ(\omega)}{d\omega} = ?$

---

$J(\omega, b)$

$\omega := \omega - \alpha \dfrac{d J(\omega, b)}{d\omega}$

$\dfrac{\partial J(\omega, b)}{\partial \omega}$

$\partial$

"partial derivative"

$b := b - \alpha \dfrac{d J(\omega, b)}{db}$

$\dfrac{\partial J(\omega, b)}{\partial b}$

$\partial$

$J$

$dw$

$db$

Andrew Ng

# Intuition about derivatives



$f(a) = 3a$

$a = 2$     $f(a) = 6$

$a = 2.001$     $f(a) = 6.003$

$\dfrac{0.003}{0.001}$   $\dfrac{\text{height}}{\text{width}}$ → slope (derivative) of $f(a)$

at $a = 2$ is 3

$a = 5$     $f(a) = 15$

$a = 5.001$     $f(a) = 15.003$

slope at $a = 5$ is also 3

$\dfrac{d\,f(a)}{da} = 3 = \dfrac{d}{da} f(a)$

0.001
0.00000001
0.00000000001


Andrew Ng

# Basics of Neural Network Programming

## More derivatives examples

deeplearning.ai

# Intuition about derivatives



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

$$\frac{d}{da}a^2 = 2a$$

$$0.001$$

$$(2a) \times 0.001$$

0.001

0.00000....01

$a = 2$      $f(a) = 4$

$a = 2.001$      $f(a) \approx 4.004$

$(4.004\boxed{001})$

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   $\boxed{a = 5}$

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

Andrew Ng

# More derivative examples

$f(a) = a^2$

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

$a = 2 \qquad f(a) = 4$

$a = 2.001 \qquad f(a) \approx 4.004$

$f(a) = a^3$

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

$a = 2 \qquad f(a) = 8$

$a = 2.001 \qquad f(a) \approx 8.012$

$f(a) = \log_e(a)$

$\ln(a)$

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

$a = 2 \qquad f(a) \approx 0.69315$

$a = 2.001 \qquad f(a) \approx 0.69365$

$\ln(a)$ 0.0005

0.001

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

0.0005

0.0005

# Basics of Neural Network Programming

## Computation Graph

deeplearning.ai

# Computation Graph

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

$$\underbrace{\underbrace{\quad}_{u}}_{J}$$

$$\underbrace{\qquad\qquad}_{v}$$

$u = bc$

$V = a + u$

$J = 3v$

$a = 5$

$b = 3$

| $u = bc$ | 6 |

$V = a + u$  11

$J = 3v$  33

$c = 2$

# Basics of Neural Network Programming

## Derivatives with a Computation Graph

deeplearning.ai

# Computing derivatives

$$\frac{dJ}{da} \quad \text{``}da\text{''} = 3$$

$a = 5$

$6$

$11$

$33$

$b = 3$

$\boxed{u=bc}$

$\boxed{v = a + u}$

$\boxed{J = 3v}$

$c = 2$

$$\frac{dJ}{dv} \quad \text{``}dv\text{''} = 3$$

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

Chain Rule

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

$$3 \times 1$$

$$\boxed{\frac{dv}{da} = 1}$$

$a \rightarrow v \rightarrow J$

$J = 3v$

$v = 11 \rightarrow 11.001$

$J = 33 \rightarrow 33.003$

$a = 5 \rightarrow 5.001$

$\rightarrow v = 11 \rightarrow 11.001$

$J = 33 \rightarrow 33.003$

$\boxed{f(a) = 3a}$

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

$\boxed{J = 3v}$

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

$$\boxed{\frac{d\,Final\,Output\,Var}{d\,var}}$$

$$\frac{dJ\,dvar}{}$$

$$\text{``}dvar\text{''}$$

Andrew Ng

# Computing derivatives

$a = 5$

$\frac{dJ}{da}$ → $\frac{da}{} = 3$

$b = 3$

$\frac{dJ}{db} = \frac{db}{} = 6$

$c = 2$

→ $\frac{dc}{} = 9$

$u = bc$  **6**

$\frac{du}{} = 3$

$v = a + u$  **11**

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

$J = 3v$  **33**

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

$\underset{3}{} \quad \underset{1}{}$

$\frac{dJ}{db} = \boxed{\frac{dJ}{du}} \cdot \frac{du}{db} = \frac{6}{}$

$\underset{\to 3}{} \quad \underset{=2}{}$

$\frac{dJ}{da} = \boxed{\frac{dJ}{du}} \cdot \frac{du}{da} = 9$

$\underset{3 \times 3}{}$

$u = 6 \to 6.001$

$v = 11 \to 11.001$

$J = 33 \to 33.003$

$b = 3 \to 3.001$

$u = b \cdot c = 6 \to 6.002$    $c = 2$

$J = 33.006$    $.006$

$v = 11.002$

$J = 3v$

# Logistic regression recap

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

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

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

Forward propagation to compute the loss function, and backwards propagation to compute the derivatives
(params w and b that minimises the loss function)

2 feature
example

$x_1$

$w_1$

$x_2$

$w_2$

$b$

$$z = w_1 x_1 + w_2 x_2 + b$$

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

$$\mathcal{L}(a, y)$$

Andrew Ng

# Logistic regression derivatives

$x_1$

$w_1$

$x_2$

$w_2$

b

$$z = w_1 x_1 + w_2 x_2 + b$$

$$a = \sigma(z)$$

$$\mathcal{L}(a, y)$$

$$dz = \frac{d\mathcal{L}}{dz} = \frac{d\mathcal{L}(a,y)}{dz}$$

$$"da" = \frac{d\mathcal{L}(a,y)}{da}$$

$$= a - y$$

$$a(1-a)$$

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

$$= \frac{d\mathcal{L}}{da} \cdot \frac{da}{dz}$$

$$\frac{d\mathcal{L}}{dw_1} = "dw_1" = x_1 \cdot dz.$$

$$dw_2 = x_2 \cdot dz.$$

$$db = dz.$$

$$w_1 := w_1 - \alpha \, dw_1$$

$$w_2 := w_2 - \alpha \, dw_2$$

$$b := b - \alpha \, db.$$

The final goal of back propagation - calculate how much we need to change the params w and b.
The params then get updated by this calculated change (derivatives)N

Andrew Ng

Basics of Neural
Network Programming

Gradient descent
on $m$ examples

deeplearning.ai

# Logistic regression on $m$ examples

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

$$\rightarrow a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

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

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

Now we expand the calculation from 1 example to across the entire training set (Cost function)

$$\frac{\partial}{\partial w_1} J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \underbrace{\frac{\partial}{\partial w_1} \mathcal{L}(a^{(i)}, y^{(i)})}_{dw_1^{(i)} - (x^{(i)}, y^{(i)})}$$

Andrew Ng

# Logistic regression on $m$ examples

$J = 0 \; ; \; dw_1 = 0 \; ; \; dw_2 = 0 \; ; \; db = 0$

$\rightarrow$ For $i = 1$ to $m$

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

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

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

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

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

$\qquad dw_2 \mathrel{+}= x_2^{(i)} dz^{(i)}$ $\qquad n = 2$

$\qquad db \mathrel{+}= dz^{(i)}$

$dw_3$
$dw_n$

$J \mathrel{/}= m \Leftarrow$

$dw_1 \mathrel{/}= m \; ; \quad dw_2 \mathrel{/}= m \; ; \; db \mathrel{/}= m. \Leftarrow$

$dw_1 = \dfrac{\partial J}{\partial w_1}$

$w_1 := w_1 - \alpha \, dw_1$

$w_2 := w_2 - \alpha \, dw_2$

$b := b - \alpha \, db.$

Vectorization

instead of For loops, for more efficiency and speed

Andrew Ng

# Basics of Neural Network Programming

## Vectorization

# What is vectorization?

$$w \in \mathbb{R}^{n_x}$$

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

$$w = \begin{bmatrix} : \\ : \end{bmatrix} \qquad x = \begin{bmatrix} : \\ : \end{bmatrix} \qquad x \in \mathbb{R}^{n_x}$$

Non-vectorized:

$$z = 0$$
for i in range$(n-x)$:
$$z \mathrel{+}= w[i] * x[i]$$

$$z \mathrel{+}= b$$

Vectorized

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

→ GPU ⎫
→ CPU ⎭ SIMD — single instruction multiple data.

Andrew Ng

Basics of Neural
Network Programming

More vectorization
examples

deeplearning.ai

# Neural network programming guideline

Whenever possible, avoid explicit for-loops.

$$u = Av$$

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

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

$$\text{for } i \ldots$$
$$\quad \text{for } j \ldots$$
$$\qquad u[i] \mathrel{+}= A[i][j] * v[j]$$

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

# 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])
```

import numpy as np

$u = np.exp(v)$

np.log(v)

np.abs(v)

np.maximum(v,0)

$v ** 2$          $1/v$

# Logistic regression derivatives

$dw = np.zeros((n_-x, 1))$

J = 0, ~~dw1 = 0, dw2 = 0~~, db = 0

→ for i = 1 to n:

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

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

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

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

for j=1...n_x
$dw_j \mathrel{+}= ...$

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

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

$n_x = 2$

$dw \mathrel{+}= x^{(i)} dz^{(i)}$

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

eliminate the second for loop that loops through the number of features / params

J = J/m, ~~$dw_1 = dw_1/m, dw_2 = dw_2/m$~~, db = db/m

$dw \mathrel{/}= m.$

Andrew Ng

# Vectorizing Logistic Regression

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

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

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

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

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

$$w^T \begin{bmatrix} x^{(1)} & x^{(2)} & \cdots & x^{(m)} \end{bmatrix}$$

$$Z = \begin{bmatrix} z^{(1)} & z^{(2)} & \cdots & z^{(m)} \end{bmatrix} = w^T X + \begin{bmatrix} b & b & \cdots & b \end{bmatrix} = \begin{bmatrix} w^T x^{(1)} + b & w^T x^{(2)} + b & \cdots & w^T x^{(m)} + b \end{bmatrix}$$

$$1 \times m$$

$$1 \times m$$

$$Z = np.dot(w.T, X) + b$$

$$(1,1) \quad \mathbb{R}$$

"Broadcasting"

$$A = \begin{bmatrix} a^{(1)} & a^{(2)} & \cdots & a^{(m)} \end{bmatrix} = \sigma(Z)$$

Andrew Ng

# Basics of Neural Network Programming

deeplearning.ai

## Vectorizing Logistic Regression's Gradient Computation

# Vectorizing Logistic Regression

$$dz^{(1)} = a^{(1)} - y^{(1)} \qquad dz^{(2)} = a^{(2)} - y^{(2)} \qquad \cdots$$

$$\boxed{dz} = [dz^{(1)} \; dz^{(2)} \cdots dz^{(m)}] \quad \Leftarrow$$
$$\underset{1 \times m}{}$$

$$A = [a^{(1)} \cdots a^{(m)}]. \qquad Y = [y^{(1)} \cdots y^{(m)}]$$

$$\rightarrow dz = A - Y = [a^{(1)} - y^{(1)} \quad a^{(2)} - y^{(2)} \quad \cdots]$$

$$\rightarrow dw = 0$$
$$dw \mathrel{+}= x^{(1)} dz^{(1)}$$
$$dw \mathrel{+}= x^{(2)} dz^{(2)}$$
$$\vdots$$
$$dw / = m$$

$$db = 0$$
$$db \mathrel{+}= dz^{(1)}$$
$$db \mathrel{+}= dz^{(2)}$$
$$\vdots$$
$$db \mathrel{+}= dz^{(m)}$$
$$db / = m.$$

$$db = \frac{1}{m} \sum_{i=1}^{m} dz^{(i)}$$

$$= \frac{1}{m} \, np.sum(dz)$$

$$dw = \frac{1}{m} X \, dz^{T}$$

$$= \frac{1}{m} \begin{bmatrix} x^{(1)} \cdots x^{(m)} \\ 1 \quad\quad 1 \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} \left[ x^{(1)} dz^{(1)} + \cdots + x^{(m)} dz^{(m)} \right]$$
$$n \times 1$$

Andrew Ng

# Implementing Logistic Regression

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

for i = 1 to m:

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

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

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

$\qquad dz^{(i)} = a^{(i)} - y^{(i)} \quad \Leftarrow$

$\qquad \left.\begin{array}{l} dw_1 \mathrel{+}= x_1^{(i)} dz^{(i)} \\ dw_2 \mathrel{+}= x_2^{(i)} dz^{(i)} \end{array}\right\} \quad dw \mathrel{+}= x^{(i)} * dz^{(i)}$

$\qquad db \mathrel{+}= dz^{(i)}$

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

$db = db/m$

for iter in range(1000): $\Leftarrow$

$\qquad Z = w^T X + b$

$\qquad\qquad = np.dot(w.T, X) + b$

$\qquad A = \sigma(Z)$

$\qquad dZ = A - Y$

$\qquad dw = \frac{1}{m} X \, dZ^T$

$\qquad db = \frac{1}{m} np.sum(dZ)$

$\qquad w := w - \alpha \, dw$

$\qquad b := b - \alpha \, db$

Andrew Ng

# Basics of Neural Network Programming

# Broadcasting in Python

deeplearning.ai

# Broadcasting example

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

$$
\begin{array}{cccc}
 & \text{Apples} & \text{Beef} & \text{Eggs} & \text{Potatoes} \\
\text{Carb} & 56.0 & 0.0 & 4.4 & 68.0 \\
\text{Protein} & 1.2 & 104.0 & 52.0 & 8.0 \\
\text{Fat} & 1.8 & 135.0 & 99.0 & 0.9
\end{array} = A
$$

$(3,4)$

$59\, cal$

$\dfrac{56}{59} \approx 94.9\%$

Calculate % of calories from Carb, Protein, Fat. Can you do this without explicit for-loop?

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

$\uparrow (3,4) \quad / \quad (1,4)$

# Broadcasting example

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \quad + \quad \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} \quad 100$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad + \quad \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$$

$(m,n) \quad (2,3) \qquad (1,n) \rightsquigarrow (m,n) \qquad (2,3)$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad + \quad \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} \quad =$$

$(m,n) \qquad\qquad (m,1)$
$\qquad\qquad\qquad \xi$
$\qquad\qquad (m,n)$

# General Principle

$$(m, n)$$

matrix

$$\begin{aligned} &+ \\ &- \\ &* \\ &/ \end{aligned}$$

$$(1, n) \quad \rightsquigarrow \quad (m, n)$$

$$(m, 1) \quad \rightsquigarrow \quad (m, n)$$

$$(m, 1) \qquad + \qquad \mathbb{R}$$

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

$$[1 \ 2 \ 3] \qquad + \qquad 100 \qquad = \qquad [101 \quad 102 \quad 103]$$

Matlab/Octave: bsxfun

# Basics of Neural Network Programming

Explanation of logistic regression cost function (Optional)

deeplearning.ai

# Logistic regression cost function

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

Interpret $\quad \hat{y} = P(y=1 \mid x)$

If $\quad y=1 \quad : \quad P(y \mid x) = \hat{y}$

If $\quad y=0 \quad : \quad P(y \mid x) = 1 - \hat{y}$

# Logistic regression cost function

$$\text{If} \quad y = 1: \qquad p(y|x) = \hat{y}$$

$$\text{If} \quad y = 0: \qquad p(y|x) = 1 - \hat{y}$$

$$p(y|x)$$

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

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

$\qquad \qquad \qquad \qquad \qquad \qquad = 1$

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

$\qquad \qquad \qquad \qquad = 1$

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

$\qquad \qquad = -\mathcal{L}(\hat{y}, y) \downarrow$

Andrew Ng

# Cost on $m$ examples

$$\log p(\text{labels in training set}) = \log \prod_{i=1}^{m} p(y^{(i)} | x^{(i)}) \leftarrow$$

$$\log p(\,\dots\,) = \sum_{i=1}^{m} \underbrace{\log p(y^{(i)} | x^{(i)})}_{-\mathcal{L}(\hat{y}^{(i)}, y^{(i)})}$$

Maximum likelihood estimate $\nwarrow$

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

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