



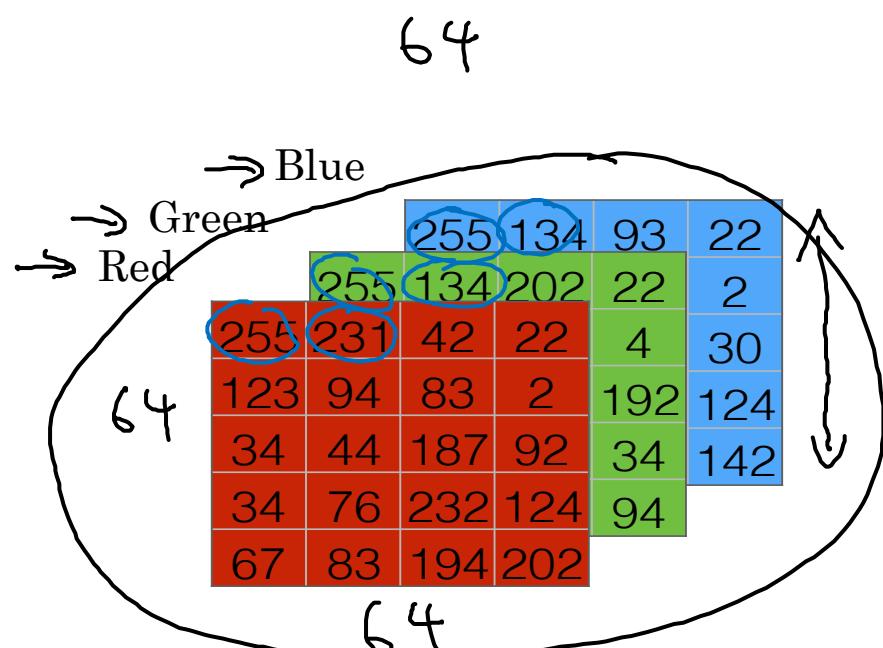
deeplearning.ai

Basics of Neural Network Programming

Binary Classification

Binary Classification

64 →  → 1 (cat) vs 0 (non cat)



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

$$64 \times 64 \times 3 = 12288$$

$$n = n_x = 12288$$

$$X \rightarrow y$$

Notation

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

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

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

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

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

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

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

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

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

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



deeplearning.ai

Basics of Neural Network Programming

Logistic Regression

Logistic Regression

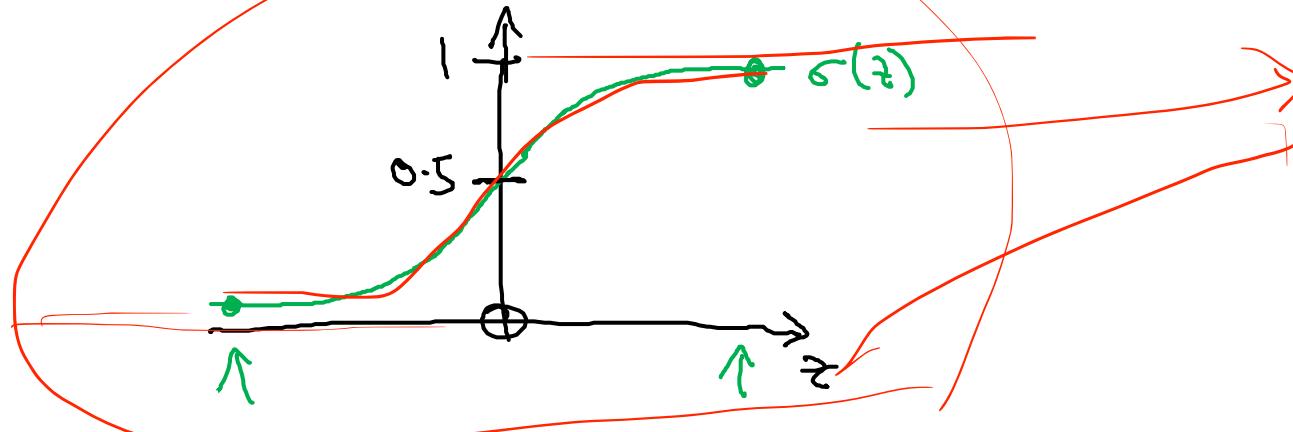
Given x , want

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

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

Output

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



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

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

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

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

$$\hat{y} = \sigma(w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots + w_{n_x} x_{n_x})$$

$w = [w_0 \ w_1 \ w_2 \ \dots \ w_{n_x}]^T$

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

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

If z large negative number

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



deeplearning.ai

Basics of Neural Network Programming

Logistic Regression cost function

Logistic Regression cost function

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

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

Loss (error) function:

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

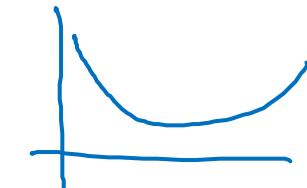
$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

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

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

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

i -th example.





deeplearning.ai

Basics of Neural Network Programming

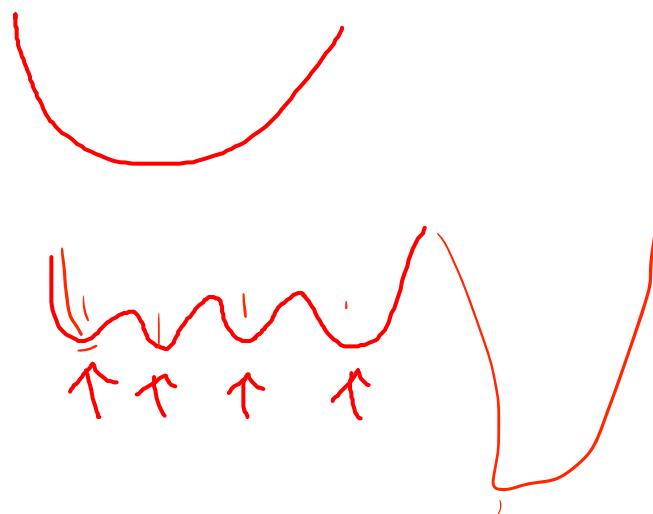
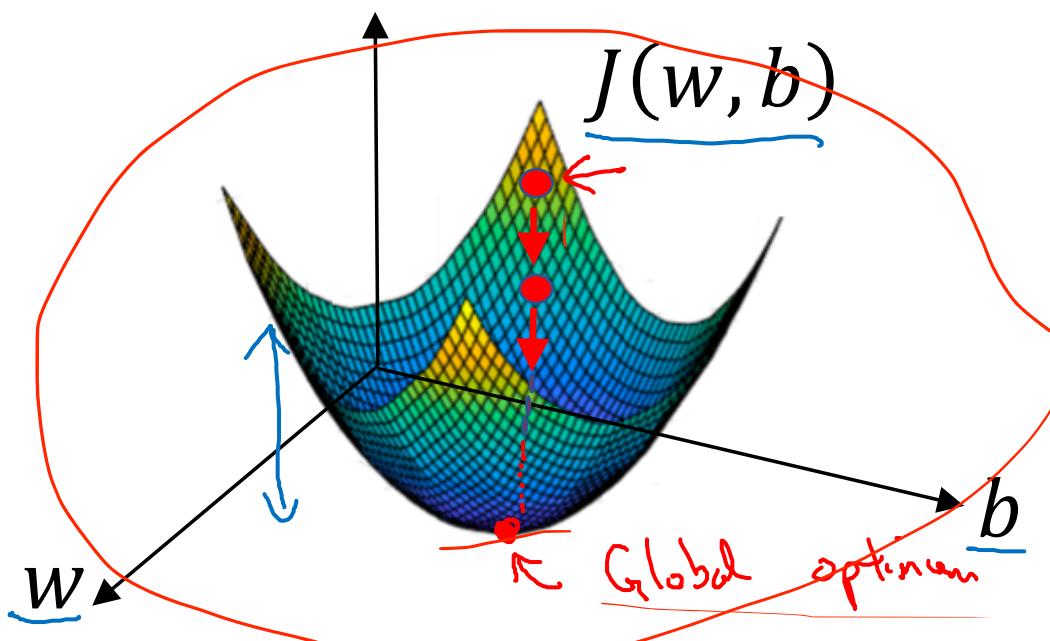
Gradient Descent

Gradient Descent

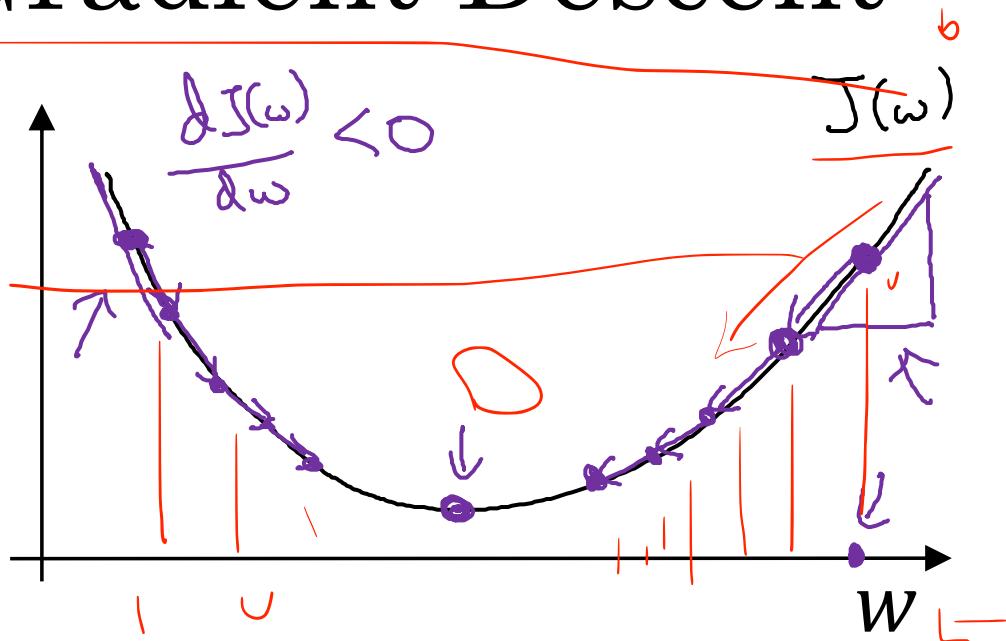
Recap: $\hat{y} = \sigma(w^T x + b)$, $\sigma(z) = \frac{1}{1+e^{-z}}$ 

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

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



Gradient Descent



Repeat {

$\omega := \omega - \frac{\alpha}{\text{learning rate}} \frac{\partial J(\omega)}{\partial \omega}$

$\omega := \omega - \alpha \frac{\partial J(\omega)}{\partial \omega}$

$\frac{\partial J(\omega)}{\partial \omega} = ?$

$$\begin{aligned} J(\omega, b) \\ \hline \omega := \omega - \alpha \frac{d J(\omega, b)}{d \omega} \\ b := b - \alpha \frac{d J(\omega, b)}{d b} \end{aligned}$$

Diagram illustrating the backpropagation of gradients for the cost function $J(\omega, b)$. The top row shows the update for ω with a red box around the gradient term $\frac{d J(\omega, b)}{d \omega}$. The bottom row shows the update for b with a red box around the gradient term $\frac{d J(\omega, b)}{d b}$. Red arrows point from these gradient terms to a red box containing the partial derivative $\frac{\partial J(\omega, b)}{\partial \omega}$ and $\frac{\partial J(\omega, b)}{\partial b}$ respectively. To the right, a blue circle labeled α is shown with a red arrow pointing to the update term $\alpha \frac{d J(\omega, b)}{d \omega}$. Another blue circle labeled d is shown with a red arrow pointing to the update term $\alpha \frac{d J(\omega, b)}{d b}$. A blue circle labeled J is shown with a red arrow pointing to the red box containing the partial derivatives. The text "partial derivative" is written in red next to the blue circle J .

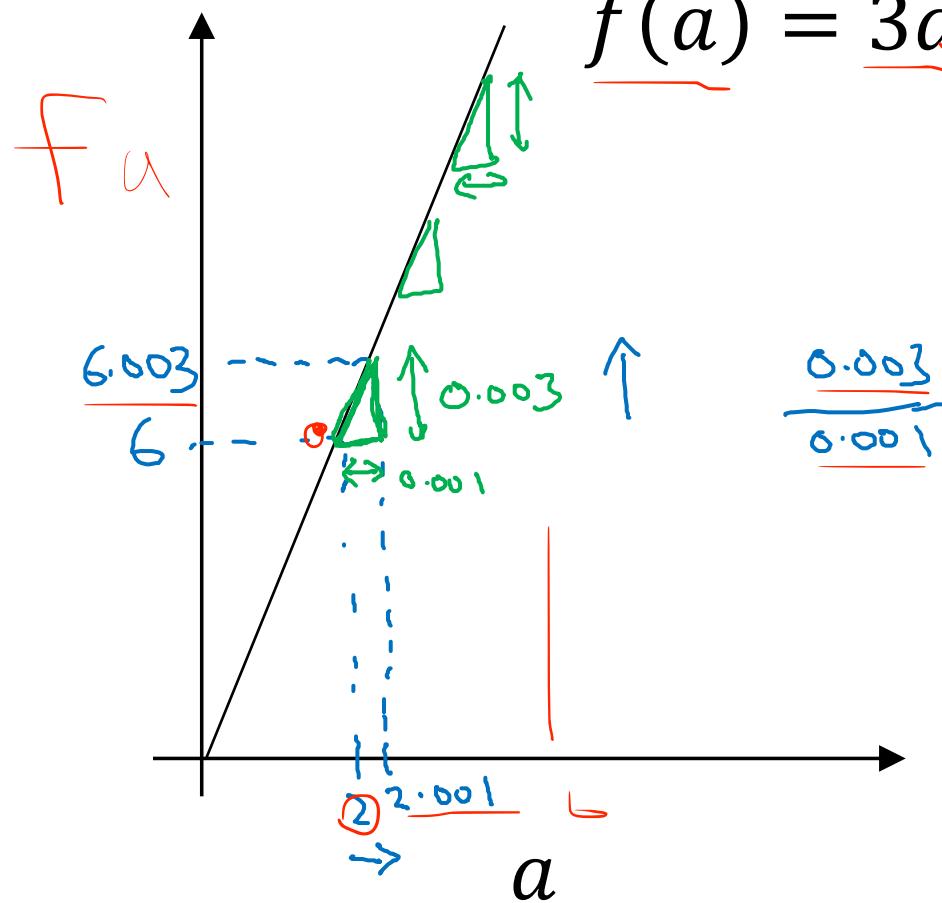


deeplearning.ai

Basics of Neural Network Programming

Derivatives

Intuition about derivatives



$$\rightarrow a = 2 \quad f(a) = 6$$
$$a = 2.001 \quad f(a) = 6.003$$

slope (derivative) of $f(a)$

at $a=2$ is 3

$$\rightarrow a = 5 \quad f(a) = 15$$
$$a = 5.001 \quad f(a) = 15.003$$

slope at $a=5$ is also 3

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

0.001 ←
0.00000001
0.0000000001

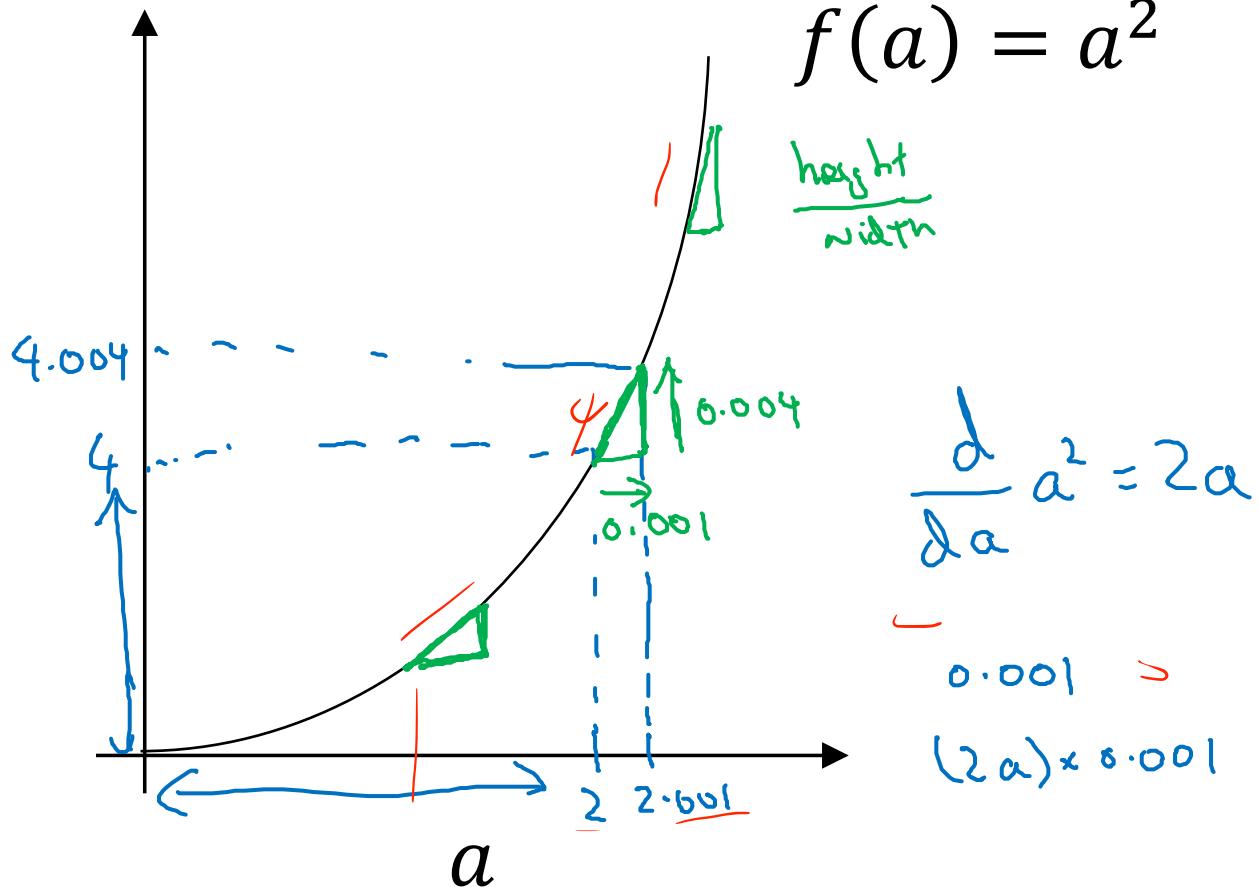


deeplearning.ai

Basics of Neural Network Programming

More derivatives
examples

Intuition about derivatives



$a = 2$

$a = 2.001$

$f(a) = 4$

$f(a) \approx 4.004$

$\frac{(4.004 - 4)}{0.001}$

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

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

$a = 5$

$a = 5.001$

$f(a) = 25$

$f(a) \approx 25.010$

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

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

More derivative examples

$$\underline{f(a) = a^2}$$

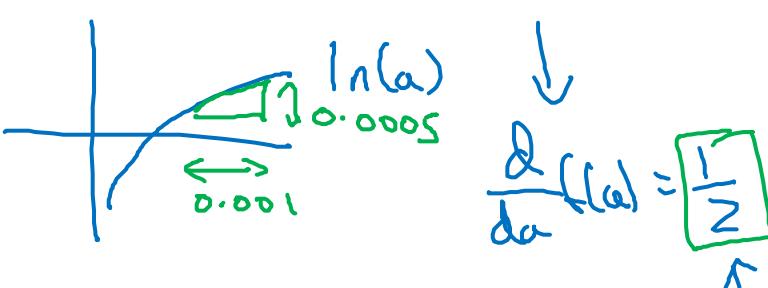
$$\frac{\partial}{\partial a} f(a) = \underbrace{2a}_{4}$$

$$\underline{f(a) = a^3}$$

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

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

$$\frac{\partial}{\partial a} f(a) = \frac{1}{a}$$



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

$$a = 2$$

$$a = 2.001$$

$$f(a) = 4$$

$$f(a) \approx 4.004$$

$$a = 2$$

$$a = \underline{2.001}$$

$$f(a) = 8$$

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

$$\downarrow a = 2$$

$$\downarrow a = \underline{2.001}$$

$$\downarrow f(a) \approx 0.69315$$

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

$$\downarrow \frac{0.0005}{0.0005}$$



deeplearning.ai

Basics of Neural Network Programming

Computation Graph

Computation Graph

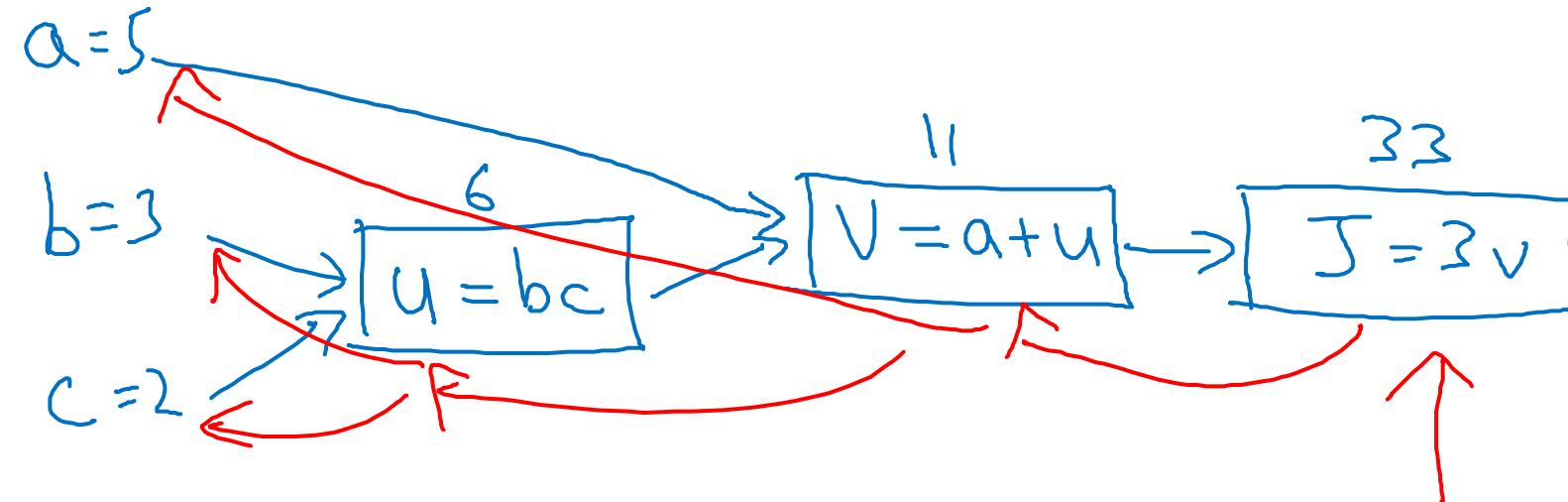
$$J(a, b, c) = 3(a + \underbrace{bc}_u) = 3(5 + 3 \times 2) = 33$$

$\underbrace{}_u$
 $\underbrace{}_v$
 $\underbrace{}_J$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



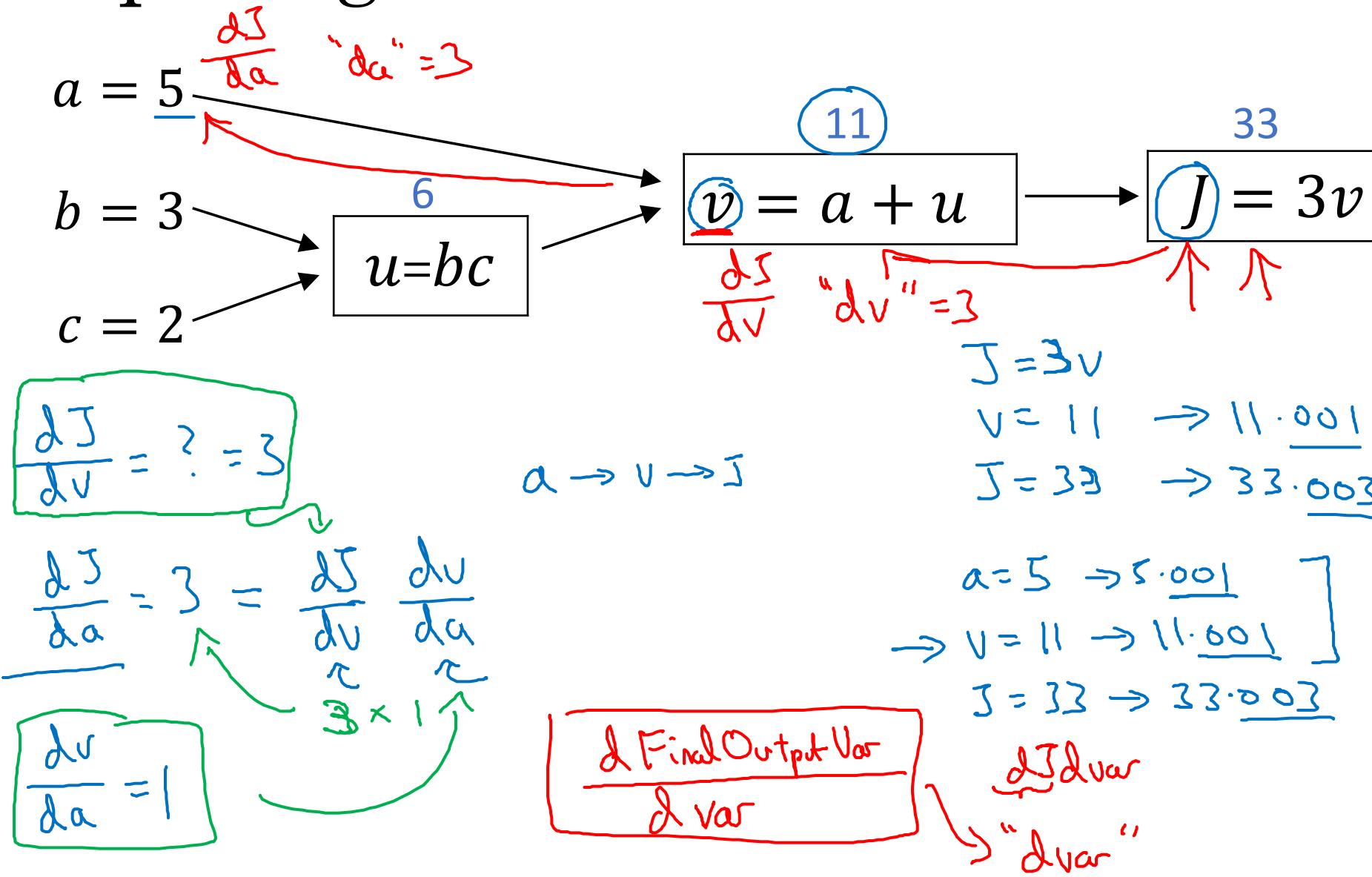


deeplearning.ai

Basics of Neural Network Programming

Derivatives with a Computation Graph

Computing derivatives



$f(a) = 3a$
 $\frac{df(b)}{da} = \frac{df}{da} = 3$
 $J = 3v$
 $\frac{dJ}{dv} = 3$

$a = 5 \rightarrow 5.001$
 $\rightarrow v = 11 \rightarrow 11.001$

$J = 33 \rightarrow 33.003$

Computing derivatives

$\frac{\partial J}{\partial a} \rightarrow \frac{\partial a}{\partial a} = 3$
 $\frac{\partial J}{\partial b} \rightarrow \frac{\partial b}{\partial b} = 6$
 $\frac{\partial J}{\partial c} \rightarrow \frac{\partial c}{\partial c} = 9$
 $\frac{\partial J}{\partial u} = 3 = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u}$
 $\frac{\partial J}{\partial b} = \boxed{3}$
 $\frac{\partial J}{\partial a} = \boxed{9}$
 $\frac{\partial J}{\partial c} = \boxed{1}$

$a = 5 \rightarrow \frac{\partial a}{\partial a} = 3$
 $b = 3 \rightarrow \frac{\partial b}{\partial b} = 6$
 $c = 2 \rightarrow \frac{\partial c}{\partial c} = 9$
 $u = bc \rightarrow \frac{\partial u}{\partial u} = 3$
 $v = a + u \rightarrow \frac{\partial v}{\partial v} = 1$
 $J = 3v \rightarrow \frac{\partial J}{\partial J} = 3$

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

$\frac{\partial J}{\partial b} = \boxed{\frac{\partial J}{\partial u}} \cdot \frac{\partial u}{\partial b} = 6$
 $\frac{\partial J}{\partial a} = \boxed{\frac{\partial J}{\partial u}} \cdot \frac{\partial u}{\partial a} = 9$

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

$v = 11.002$
 $J = 3v$



deeplearning.ai

Basics of Neural Network Programming

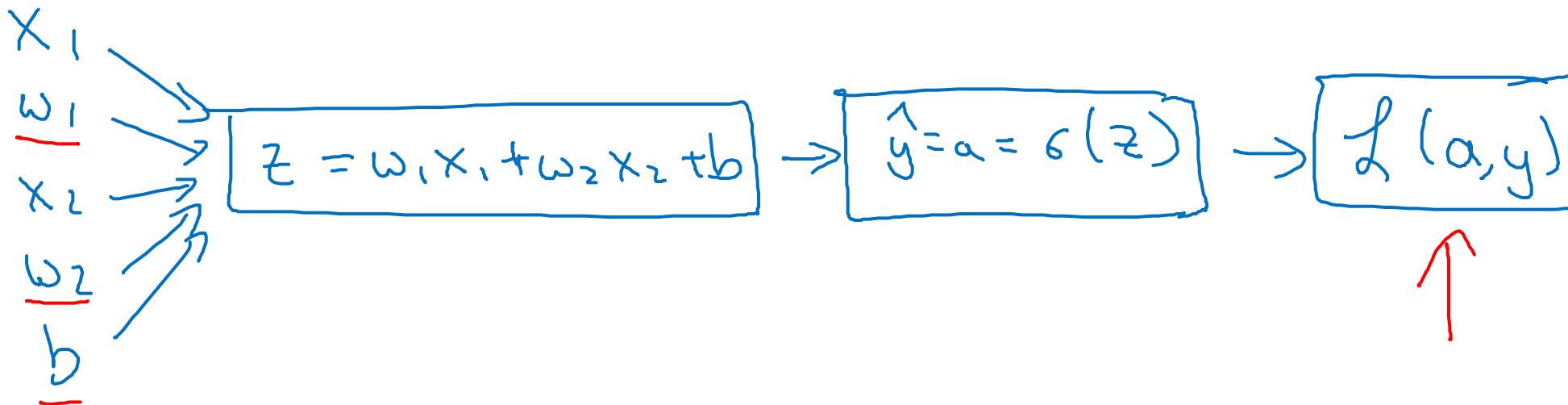
Logistic Regression Gradient descent

Logistic regression recap

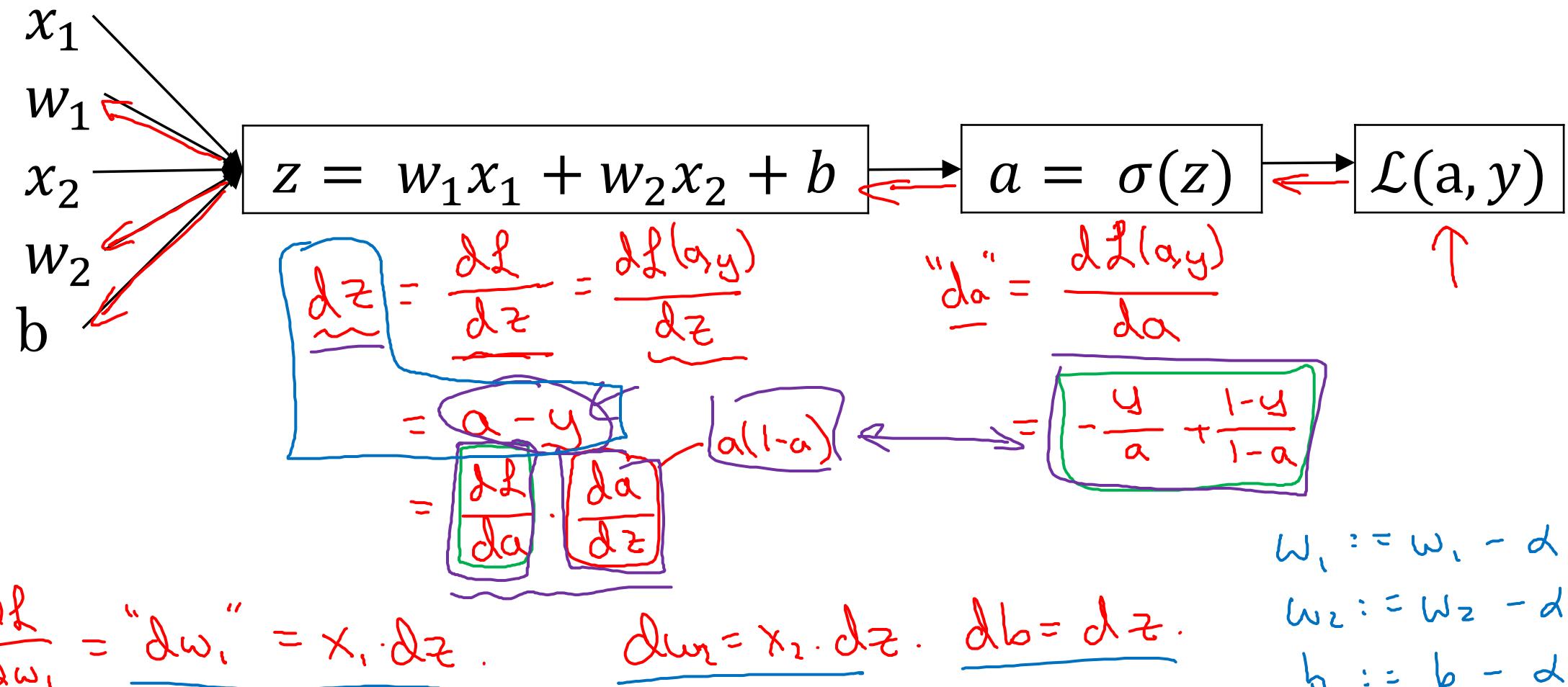
$$\rightarrow z = w^T x + b$$

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

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



Logistic regression derivatives





deeplearning.ai

Basics of Neural Network Programming

Gradient descent
on m examples

Logistic regression on m examples

$$\underline{J(\omega, b)} = \frac{1}{m} \sum_{i=1}^m \ell(a^{(i)}, y^{(i)}) \quad (x^{(i)}, y^{(i)})$$
$$\Rightarrow a^{(i)} = \hat{y}^{(i)} = g(z^{(i)}) = g(\omega^\top x^{(i)} + b) \quad \underline{dw_1^{(i)}}, \underline{dw_2^{(i)}}, \underline{db^{(i)}}$$

$$\underline{\frac{\partial}{\partial \omega_1} J(\omega, b)} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial \omega_1} \ell(a^{(i)}, y^{(i)})}_{\underline{dw_1^{(i)}}} \quad - (x^{(i)}, y^{(i)})$$

Logistic regression on m examples

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

→ For $i = 1$ to m

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

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

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

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

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

$$\begin{aligned} \Delta w_3 & \\ \vdots & \\ \Delta w_n & \end{aligned}$$

$$\uparrow \quad \uparrow \quad \uparrow \quad n=2$$

$$J / = m \leftarrow$$

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

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

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

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

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

Vectorization



deeplearning.ai

Basics of Neural Network Programming

Vectorization

What is vectorization?

$$z = \underbrace{\omega^T x}_{\text{Non-vectorized}} + b$$

Non-vectorized:

$$z = 0$$

for i in range(n - x):
 $z += \omega[i] * x[i]$

$$z += b$$

$$\omega = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

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

Vectorized

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

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



deeplearning.ai

Basics of Neural Network Programming

More vectorization examples

Neural network programming guideline

Whenever possible, avoid explicit for-loops.

$$u = Av$$

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

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

for i ...

 for j ...

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

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

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  
v/v
```

Logistic regression derivatives

$$J = 0, \boxed{dw_1 = 0, dw_2 = 0}, db = 0$$

$$dw = np.zeros((n_x, 1))$$

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 $j=1 \dots n_x$
 $dw_j += x_j^{(i)} dz^{(i)}$
 $dw_0 += dz^{(i)}$

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

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

$$dw /= m.$$



deeplearning.ai

Basics of Neural Network Programming

Vectorizing Logistic Regression

Vectorizing Logistic Regression

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

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

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

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

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

$$\underline{\underline{\omega}}^T \begin{bmatrix} 1 & x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$$

$$\underline{\underline{Z}} = \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \underline{\underline{\omega}}^T \underline{\underline{X}} + \begin{bmatrix} b & b & \dots & b \end{bmatrix}_{1 \times m} = \begin{bmatrix} \underline{\underline{\omega}}^T \underline{\underline{x}}^{(1)} + b \\ \underline{\underline{\omega}}^T \underline{\underline{x}}^{(2)} + b \\ \vdots \\ \underline{\underline{\omega}}^T \underline{\underline{x}}^{(m)} + b \end{bmatrix}$$

$$\rightarrow \underline{\underline{Z}} = \text{np.dot}(\underline{\underline{\omega}}^T, \underline{\underline{X}}) + \underline{\underline{b}}_{(1, m)}$$

R

"Broadcasting"

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



deeplearning.ai

Basics of Neural Network Programming

Vectorizing Logistic Regression's Gradient Computation

Vectorizing Logistic Regression

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)} \quad \dots$$

$$dz = \begin{bmatrix} dz^{(1)} & dz^{(2)} & \dots & dz^{(m)} \end{bmatrix} \quad \leftarrow$$

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

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

$$\begin{aligned} \rightarrow dw &= 0 \\ dw &+ = \frac{x^{(1)} dz^{(1)}}{m} \\ dw &+ = \frac{x^{(2)} dz^{(2)}}{m} \\ &\vdots \\ dw &= m \end{aligned}$$

$$\begin{aligned} db &= 0 \\ db &+ = dz^{(1)} \\ db &+ = dz^{(2)} \\ &\vdots \\ db &+ = dz^{(m)} \\ db &= m \end{aligned}$$

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

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

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

$$= \frac{1}{m} \begin{bmatrix} 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} \begin{bmatrix} x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)} \end{bmatrix}_{n \times 1}$$

Implementing Logistic Regression

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

for $i = 1$ to m :

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

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

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

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

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

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

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

$$db = db/m$$

```
for iter in range(1000):  
    z = w^T X + b  
    = np.dot(w.T, X) + b  
    A = sigma(z)  
    dZ = A - Y  
    dw = 1/m * X * dZ^T  
    db = 1/m * np.sum(dZ)  
  
    w := w - alpha * dw  
    b := b - alpha * db
```



deeplearning.ai

Basics of Neural Network Programming

Broadcasting in Python

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	
Fat	1.8	135.0	99.0	0.9	

= A
(3,4)

$$59_{\text{cal}} \quad \frac{56}{59} \approx 94.9\%$$

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

```
cal = A.sum(axis = 0)
```

$$\text{percentage} = \frac{100 * A / (\text{cal.reshape}(1, 4))}{\uparrow(3, 4) / (1, 4)}$$

Broadcasting example

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

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} \xrightarrow{(m,n) \quad (2,3) \quad (1,n) \rightsquigarrow (m,n) \quad (2,3)}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \xleftarrow{(m,n) \quad (m,1)}$$

General Principle

$$\begin{array}{ccc} (m, n) & \begin{array}{c} + \\ - \\ \times \\ \diagup \end{array} & (1, n) \rightsquigarrow (m, n) \\ \underline{\text{matrix}} & & (m, 1) \rightsquigarrow (m, n) \end{array}$$

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

Matlab/Octave: bsxfun



deeplearning.ai

Basics of Neural Network Programming

Explanation of logistic
regression cost function
(Optional)

Logistic regression cost function

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

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

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

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

Logistic regression cost function

$$\begin{aligned} \rightarrow & \boxed{\text{If } y = 1: \quad p(y|x) = \hat{y}} \\ \rightarrow & \boxed{\text{If } y = 0: \quad p(y|x) = 1 - \hat{y}} \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} p(y|x)$$
$$\boxed{p(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}} \quad \leftarrow$$
$$\text{If } y=1: \quad p(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)} = 1$$
$$\text{If } y=0: \quad p(y|x) = \hat{y}^0 (1-\hat{y})^{(1-y)} = 1 \times (1-\hat{y}) = 1 - \hat{y}$$
$$\uparrow \log p(y|x) = \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y})$$
$$= -\frac{1}{m} \sum \log (\hat{y} \cdot y) \quad \downarrow$$

Cost on m examples

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

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

Maximum likelihood
estimation \nearrow

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

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