



deeplearning.ai

Basics of Neural Network Programming

Binary Classification

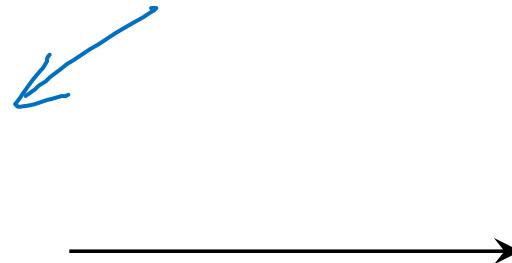
Binary Classification

Input cat image



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

Pixel intensity values for cat image



Aim to classify
1 (cat) vs 0 (non cat)

y
Input feature
Total Dimension
of the vector

Pixel intensity vector

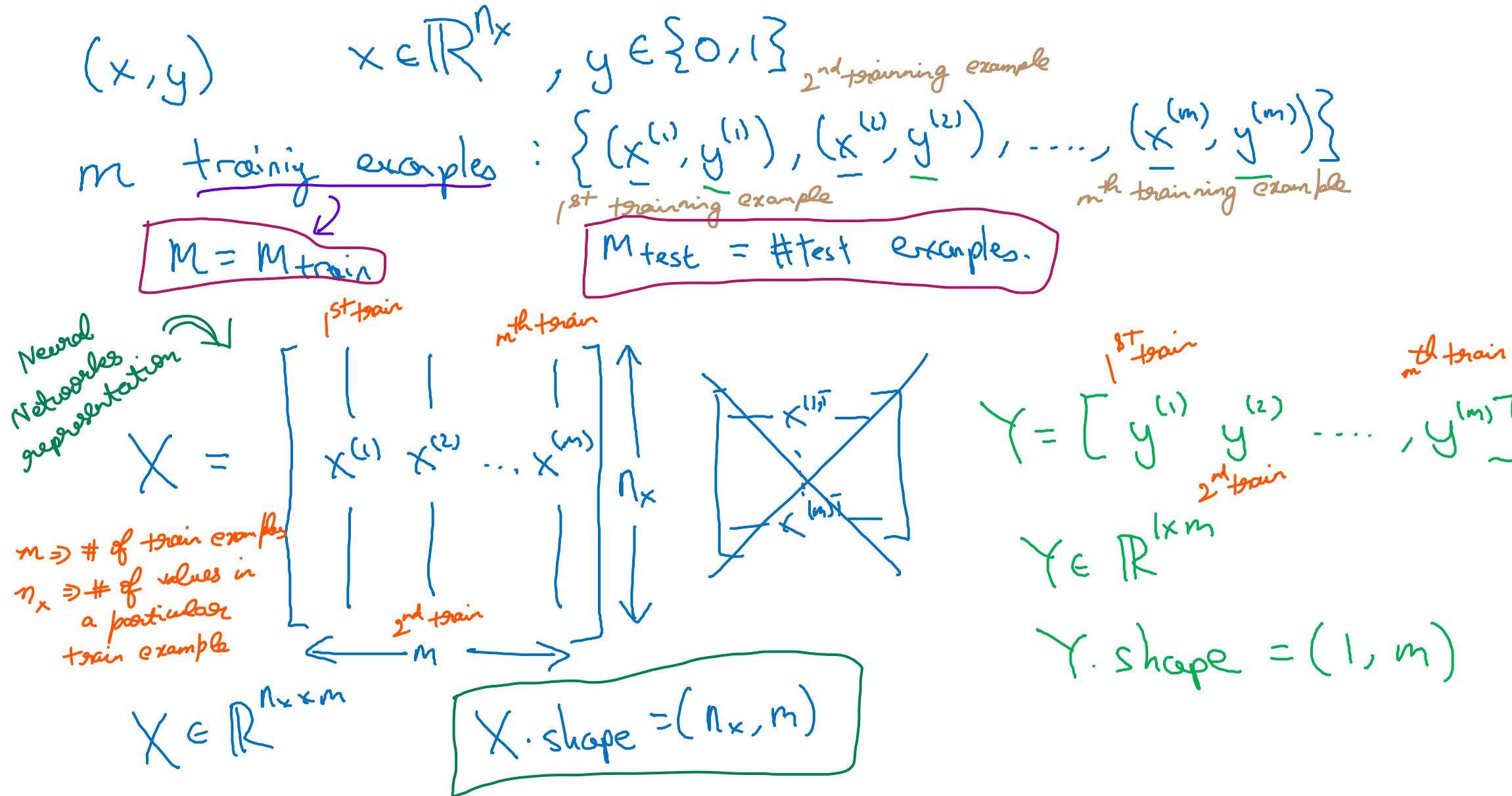
$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \\ 255 \\ 134 \end{bmatrix}$$

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

$$n = n_x = 12288$$

$$x \rightarrow y$$

Notation





deeplearning.ai

Basics of Neural Network Programming

Logistic Regression

Logistic Regression

If x is a picture, we want to tell what is the probability of a cat picture?

what is the probability of a cat picture?

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

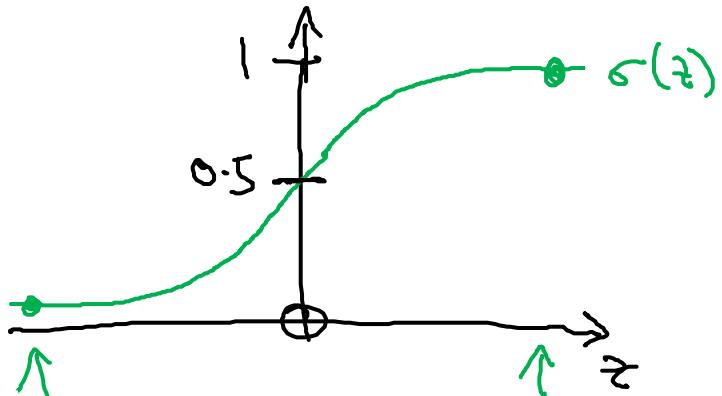
$$x_0 = 1, \quad x \in \mathbb{R}^{n_x+1}$$

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \quad \left. \begin{array}{l} \theta_0 \} b \leftarrow \text{Intercept} \\ \theta_1 \} w \\ \theta_2 \\ \vdots \\ \theta_{n_x} \} w \leftarrow \text{slope} \end{array} \right.$$

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

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



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

Not going to use this notation

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

If z large negative number

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



deeplearning.ai

Basics of Neural Network Programming

Logistic Regression cost function

Logistic Regression cost function

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

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

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

Loss (error) function:

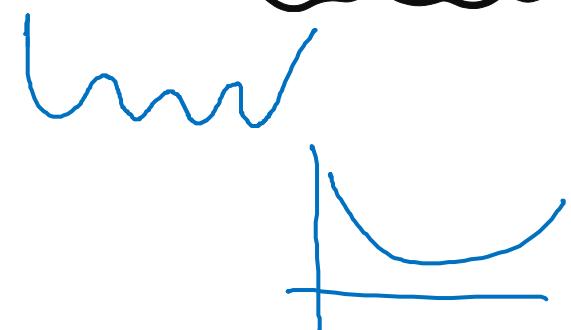
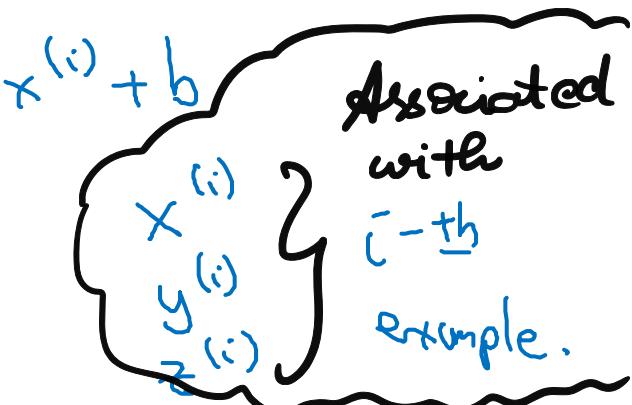
$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

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

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





deeplearning.ai

Basics of Neural Network Programming

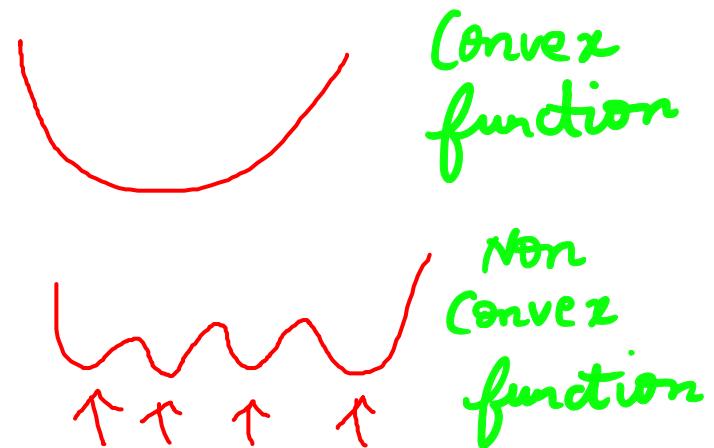
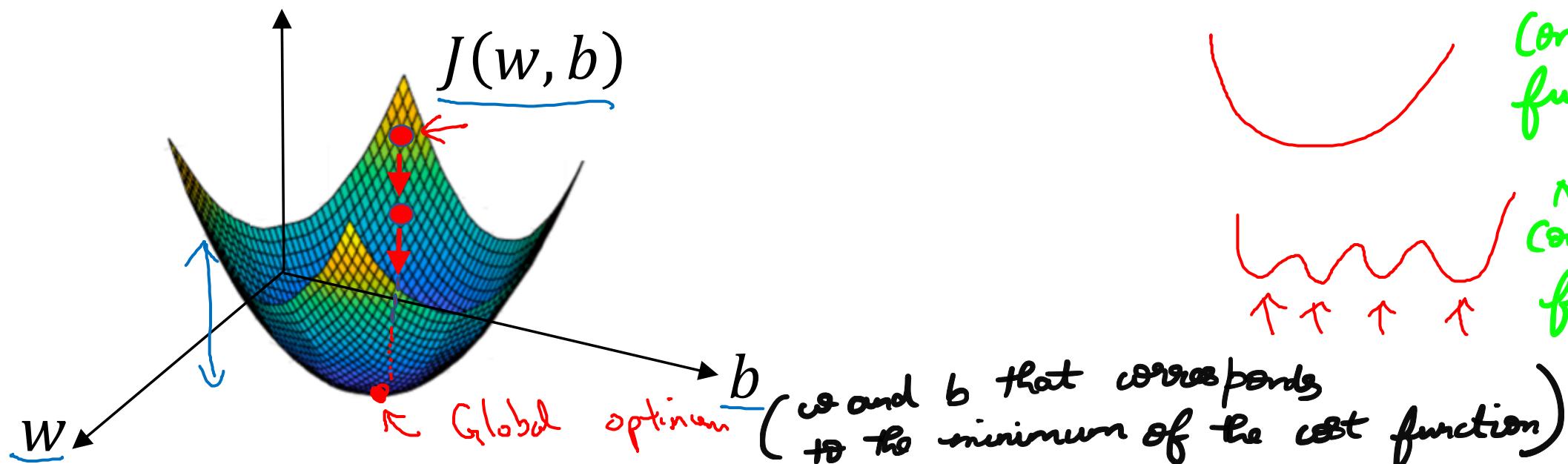
Gradient Descent

Gradient Descent

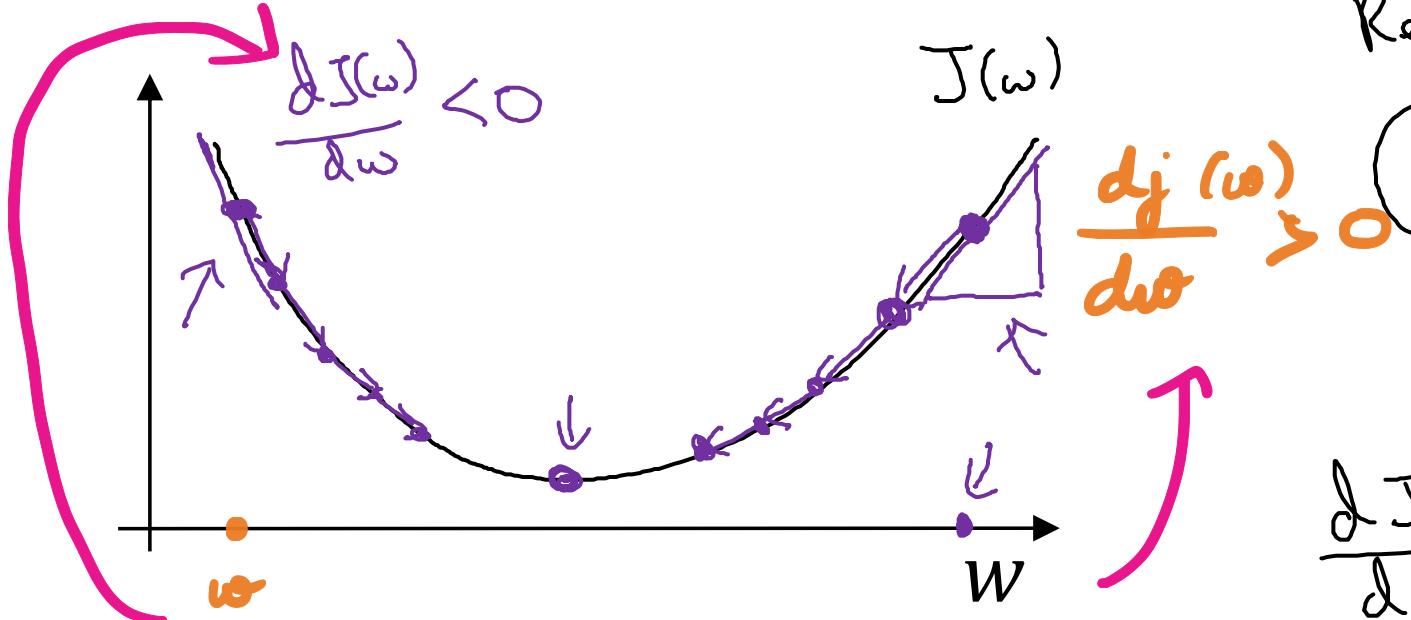
Recap: $\hat{y} = \sigma(w^T x + b)$, $\sigma(z) = \frac{1}{1+e^{-z}}$

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

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



Gradient Descent



Repeat {

$$w := w - \alpha$$

learning rate

$$\boxed{\frac{dJ(w)}{dw}}$$

$$"dw"$$

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

$$\boxed{\frac{dJ(w)}{dw}} = ?$$

$$J(w, b)$$

$$w := w - \alpha \boxed{\frac{dJ(w, b)}{dw}}$$

$$\boxed{\frac{\partial J(w, b)}{\partial w}}$$

"partial derivative"
J

$$b := b - \alpha \boxed{\frac{dJ(w, b)}{db}}$$

$$\boxed{\frac{\partial J(w, b)}{\partial b}}$$

$$\frac{dw}{db}$$

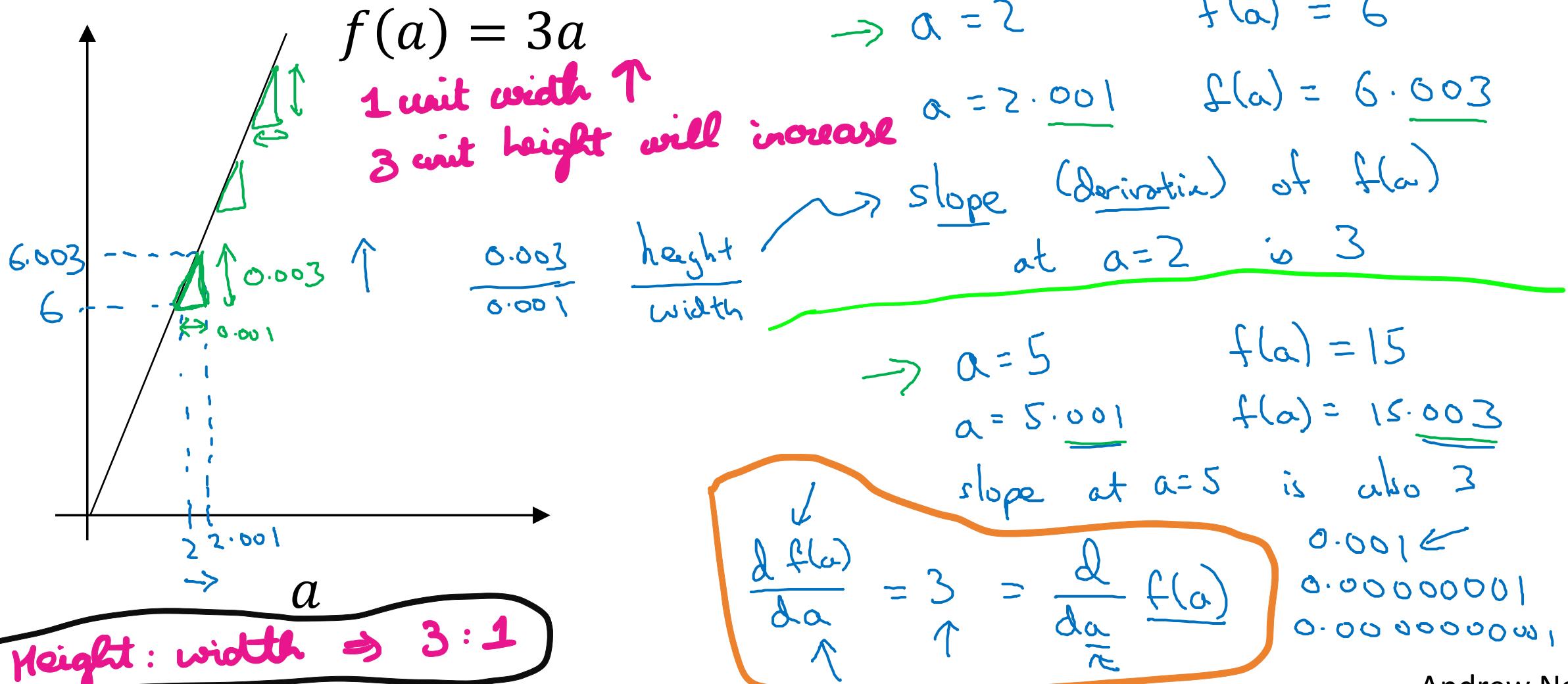


deeplearning.ai

Basics of Neural Network Programming

Derivatives

Intuition about derivatives



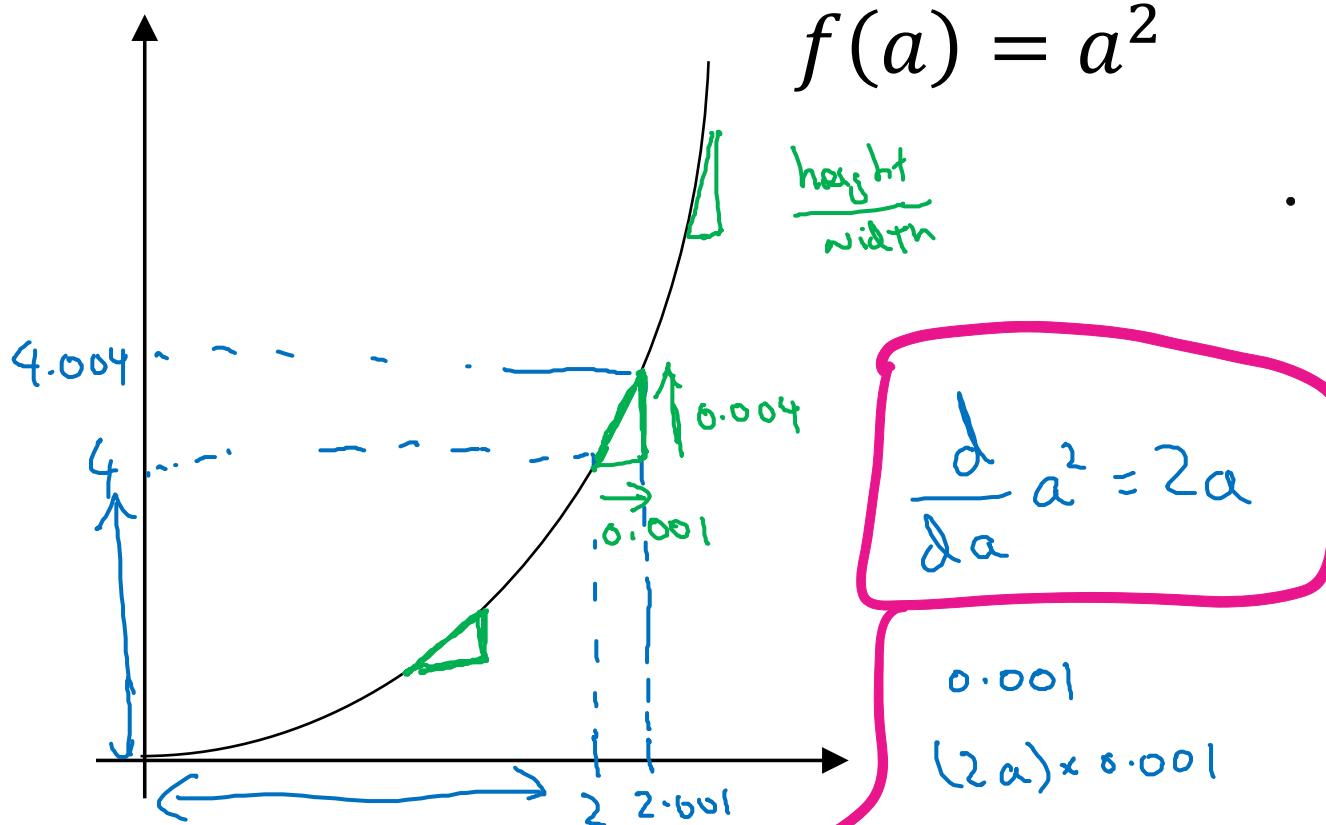


deeplearning.ai

Basics of Neural Network Programming

More derivatives
examples

Intuition about derivatives



If width increase by 1 unit, then the height increases by 2 units.

$$a=2$$

$$a=2.001$$

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

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

$$a=5$$

$$a=5.001$$

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

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

$$f(a) \approx 4.004$$

$$(4.004 \boxed{004})$$

$$\text{when } a=2$$

$$f(a) = 25$$

$$f(a) \approx 25.010$$

$$\text{when } a=5$$

More derivative examples

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

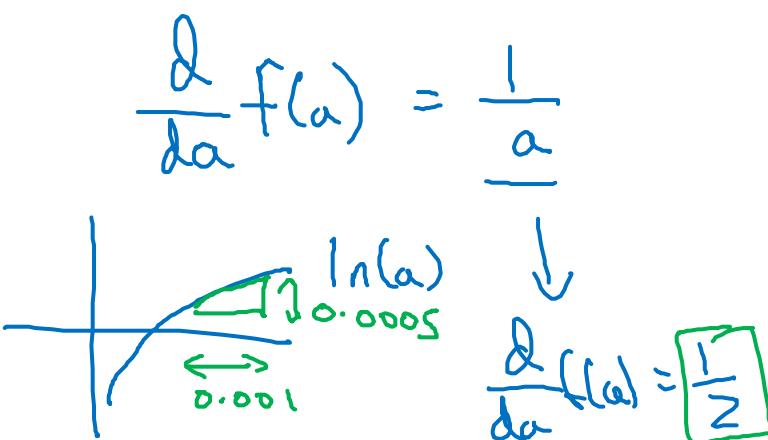
$$\frac{\partial}{\partial a} f(a) = \underbrace{2a}_{4}$$

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

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

$$f(a) = \begin{matrix} \log_e(a) \\ \ln(a) \end{matrix}$$

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



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

$$a = 2$$

$$f(a) = 4$$

$$a = 2.001$$

$$f(a) \approx 4.004$$

$$a = 2$$

$$f(a) = 8$$

$$a = \underline{2.001}$$

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

$$a = 2$$

$$f(a) \approx 0.69315$$

$$a = \underline{2.001}$$

$$f(a) \approx 0.69365$$

$$\downarrow$$

$$0.0005$$

$$\xrightarrow{0.0005}$$

$a \uparrow$ by 0.001
 $f(a) \uparrow$ by 0.0005

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

Andrew Ng



deeplearning.ai

Basics of Neural Network Programming

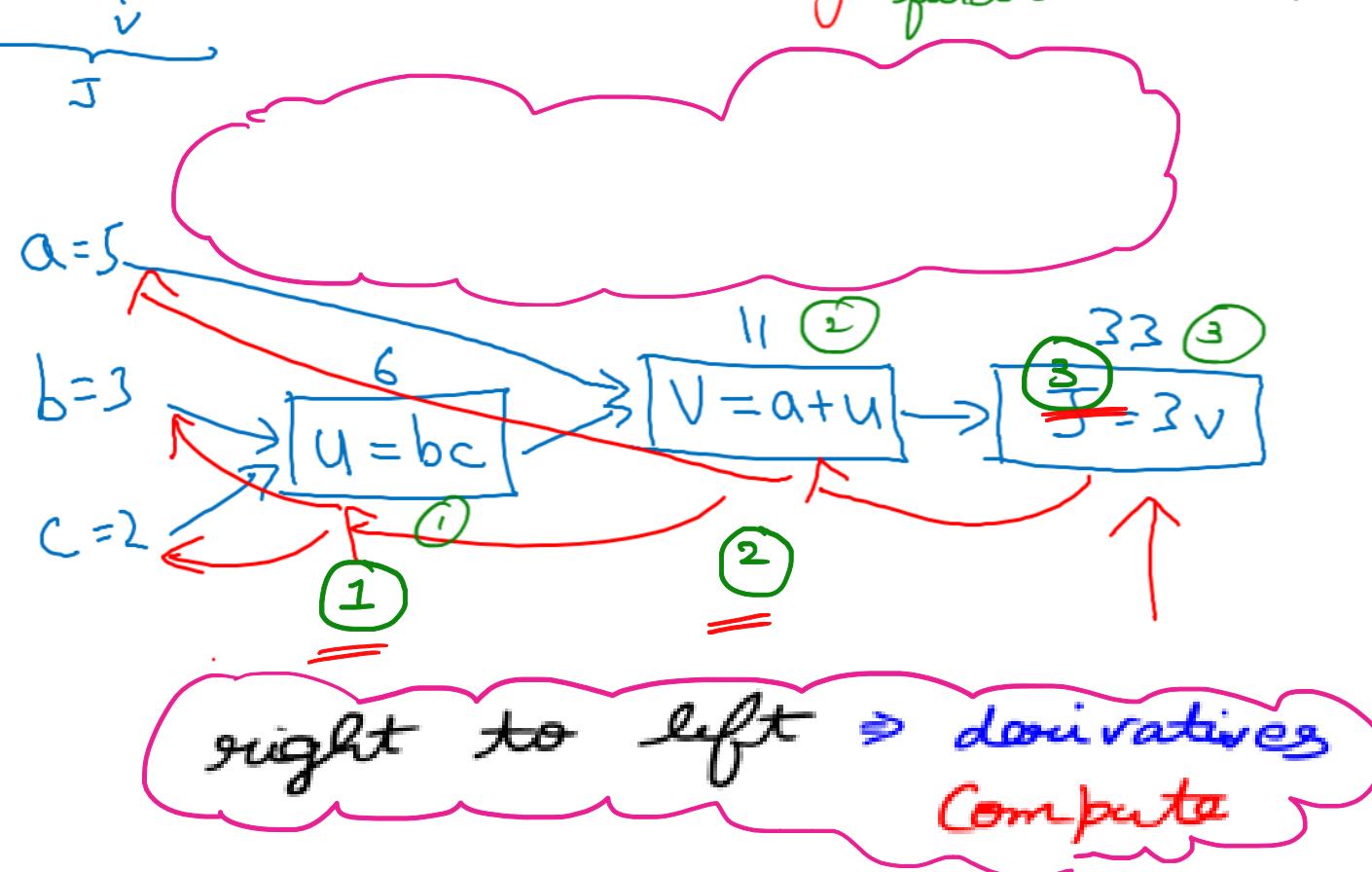
Computation Graph

Computation Graph

$$J(a, b, c) = 3(\underbrace{a + bc}_u) = 3(5 + 3 \times 2) = 23 \quad \left. \begin{array}{l} \text{Compute this} \\ \text{function } J(a, b, c) \end{array} \right\}$$

- ① $u = bc$
- ② $v = a + u$
- ③ $J = 3v$

We have to find all the three in-order to compute ① ② and ③



Computation Graph

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

$\underbrace{u}_{\downarrow}$
 $\underbrace{v}_{\downarrow}$
 $\underbrace{J}_{\downarrow}$

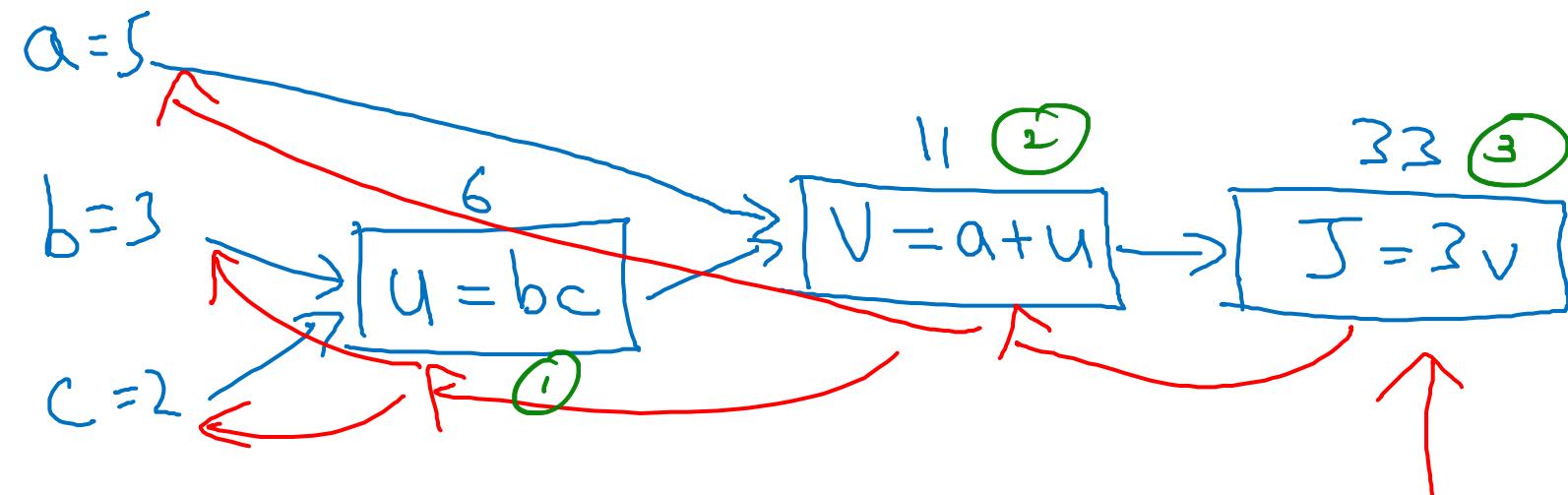
} compute this function $J(a, b, c)$

① $u = bc$

② $v = a + u$

③ $J = 3v$

We have to find all the three in-order to compute ① ② and ③



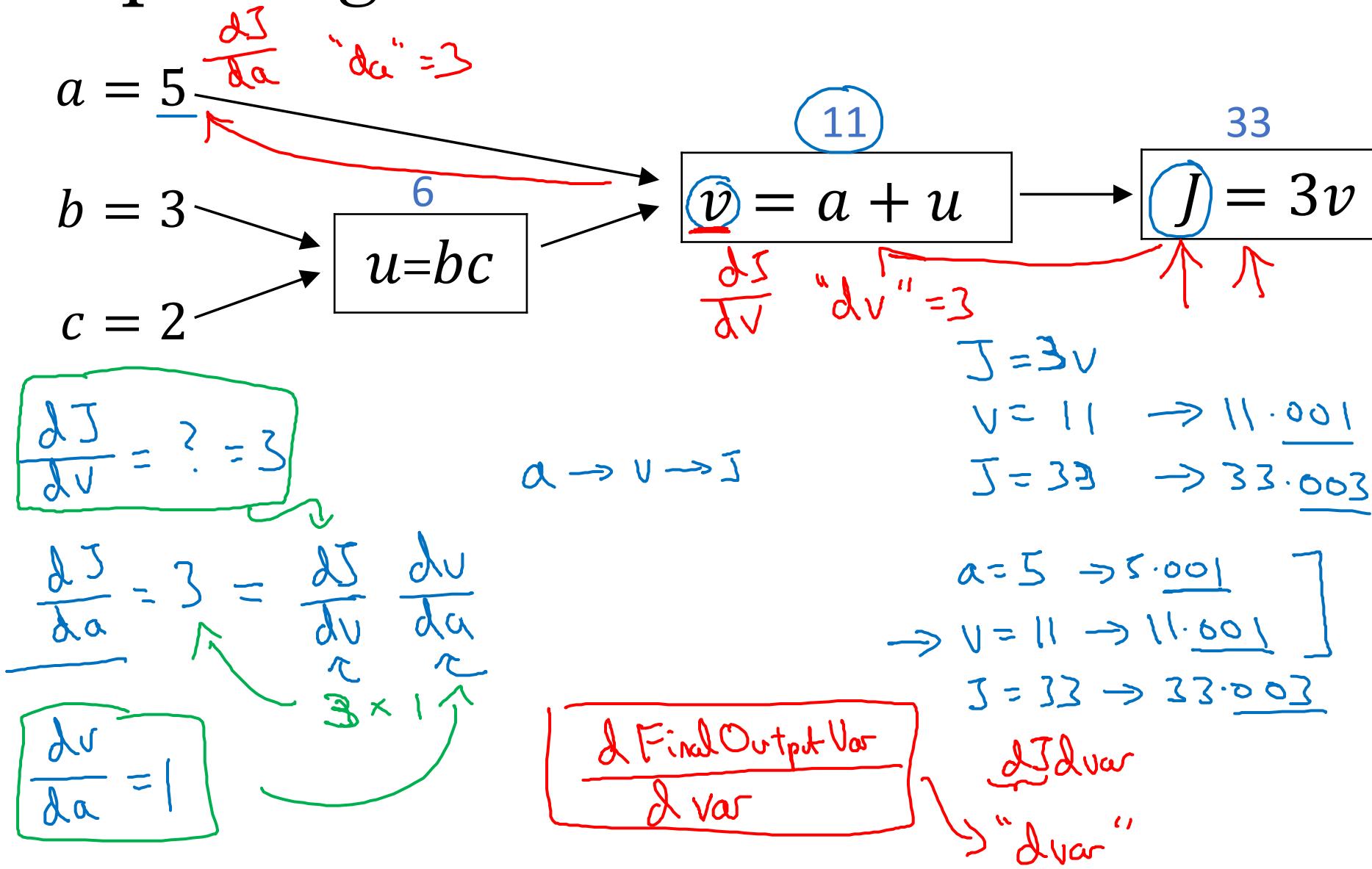


deeplearning.ai

Basics of Neural Network Programming

Derivatives with a Computation Graph

Computing derivatives



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

Computing derivatives

$\frac{\partial J}{\partial a} \rightarrow \underline{a = 5} \quad \frac{\partial a}{\partial a} = 3$
 $\frac{\partial J}{\partial b} \rightarrow \underline{b = 3} \quad \frac{\partial b}{\partial b} = 6$
 $\frac{\partial J}{\partial c} \rightarrow \underline{c = 2} \quad \frac{\partial c}{\partial c} = 9$
 $\frac{\partial J}{\partial u} = 3 = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u}$
 $\qquad\qquad\qquad \underbrace{3}_{3} \qquad\qquad\qquad \underbrace{-1}_{1}$

$a = 5 \rightarrow \underline{a = 5}$
 $b = 3 \rightarrow \underline{b = 3}$
 $c = 2 \rightarrow \underline{c = 2}$
 $u = bc \rightarrow \underline{u = 6}$
 $v = a + u \rightarrow \underline{v = 11}$
 $J = 3v \rightarrow \underline{J = 33}$

$\frac{\partial v}{\partial u} = 3 \quad \frac{\partial J}{\partial v} = 3$
 \uparrow

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

$\frac{\partial J}{\partial b} = \boxed{\frac{\partial J}{\partial u}} \cdot \underbrace{\frac{\partial u}{\partial b}}_{=2} = 6$
 $\qquad\qquad\qquad \underbrace{3}_{3}$

$\frac{\partial J}{\partial a} = \boxed{\frac{\partial J}{\partial u}} \cdot \underbrace{\frac{\partial u}{\partial a}}_{=3} = 9$
 $\qquad\qquad\qquad \underbrace{3}_{3} \times \underbrace{3}_{3}$

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

$c = 2$
 $.006$

$v = 11.002$
 $J = 3v$

$$a = 5$$

$$b = 3$$

$$c = 2$$

$$\mu = bc$$

$$v = a + \mu$$

$$J = 3v$$

$$J = 33$$

$$v = 11$$

$$J = 33 \rightarrow 33.003$$

$$v = 11.001$$

Step-1 (a 's effect)

3 times increase

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

Therefore slight change in v will increase the J , 3 times.

$$v = a + \mu$$

$$a = 5 \rightarrow 5.001 \quad \uparrow v = a + \mu$$

$$v \Rightarrow 6 + 5 \Rightarrow 11 \rightarrow 6 + 5.001 \Rightarrow 11.001 \quad \uparrow J = 3v \uparrow$$

$$J \Rightarrow 3 * 11 \Rightarrow 33 \quad 3 * 11.001 \Rightarrow 33.003$$

$$\frac{dJ}{da} \rightarrow \frac{dv}{da} * \frac{dJ}{dv}$$

1
3

$$\frac{dJ}{da} \Rightarrow 1 * 3 \Rightarrow 3$$

Step-2 (b's effect)



deeplearning.ai

Basics of Neural Network Programming

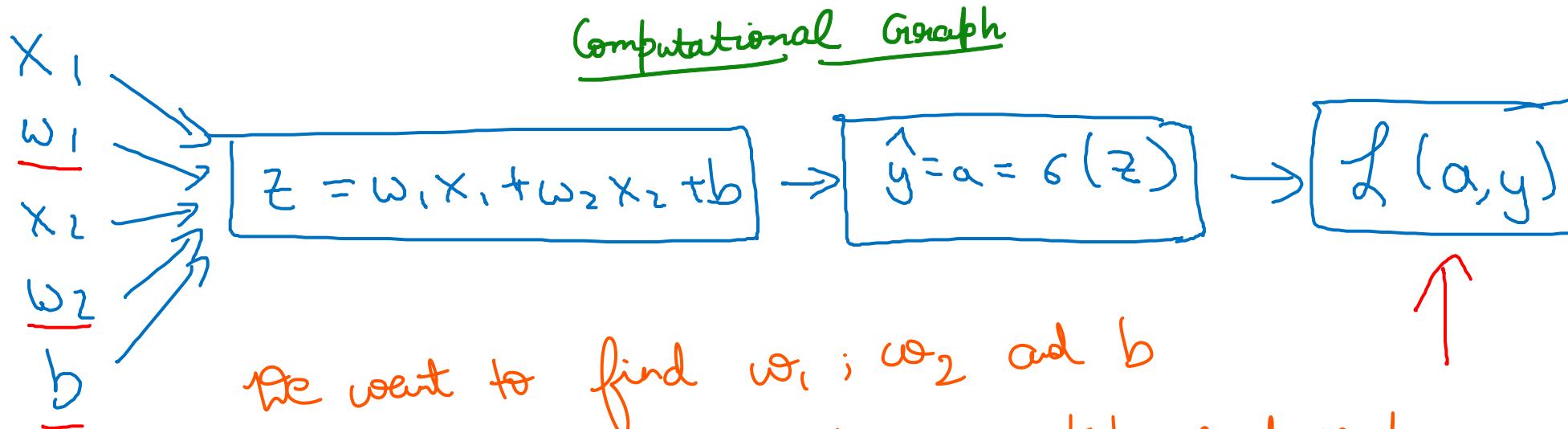
Logistic Regression Gradient descent

Logistic regression recap

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

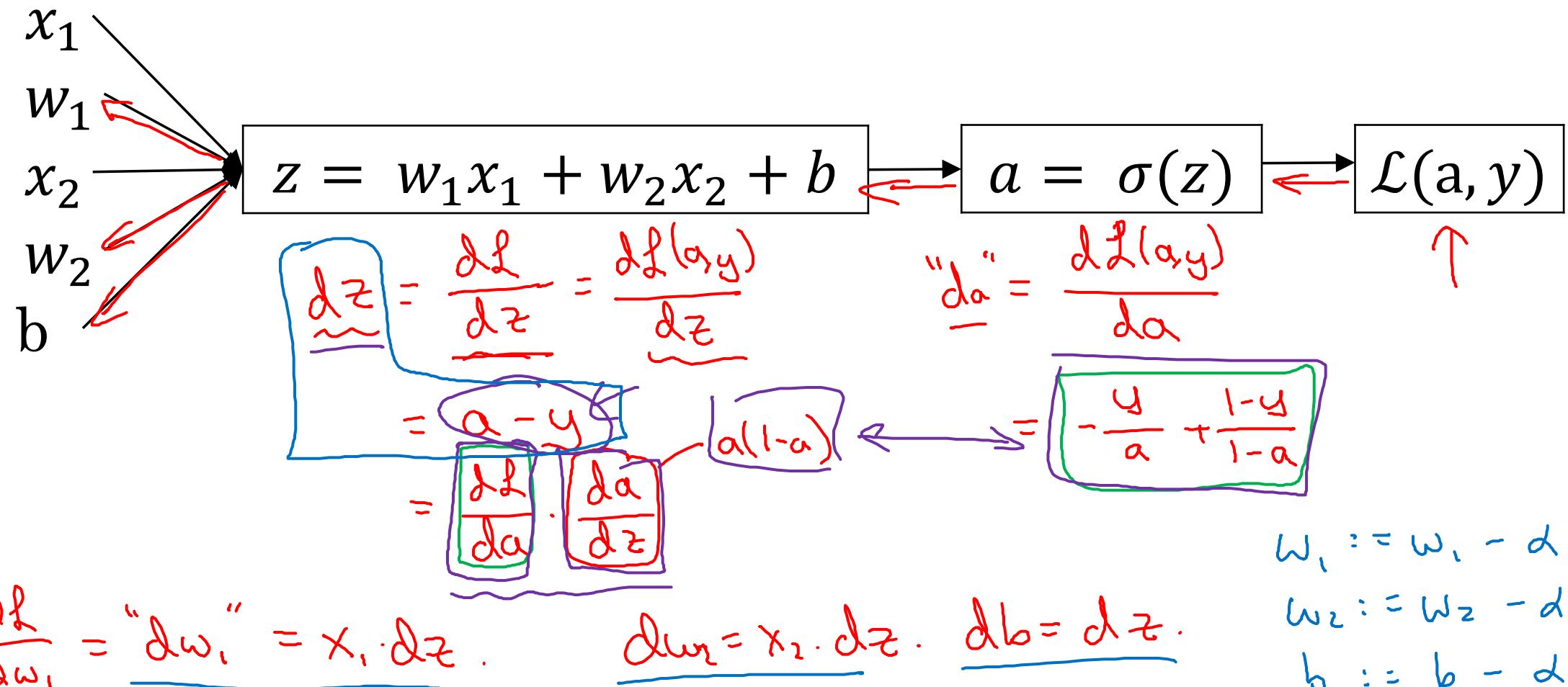
$$\rightarrow \hat{y} = a = \sigma(z)$$

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



We want to find w_1 ; w_2 and b
in order that the loss function must be reduced

Logistic regression derivatives



$$\frac{\partial \mathcal{L}}{\partial w_1} = "dw_1" = x_1 \cdot dz$$

$$dw_1 = x_1 \cdot dz. \quad db = dz.$$

$$\begin{aligned} w_1 &:= w_1 - \alpha dw_1 \\ w_2 &:= w_2 - \alpha dw_2 \\ b &:= b - \alpha db. \end{aligned}$$

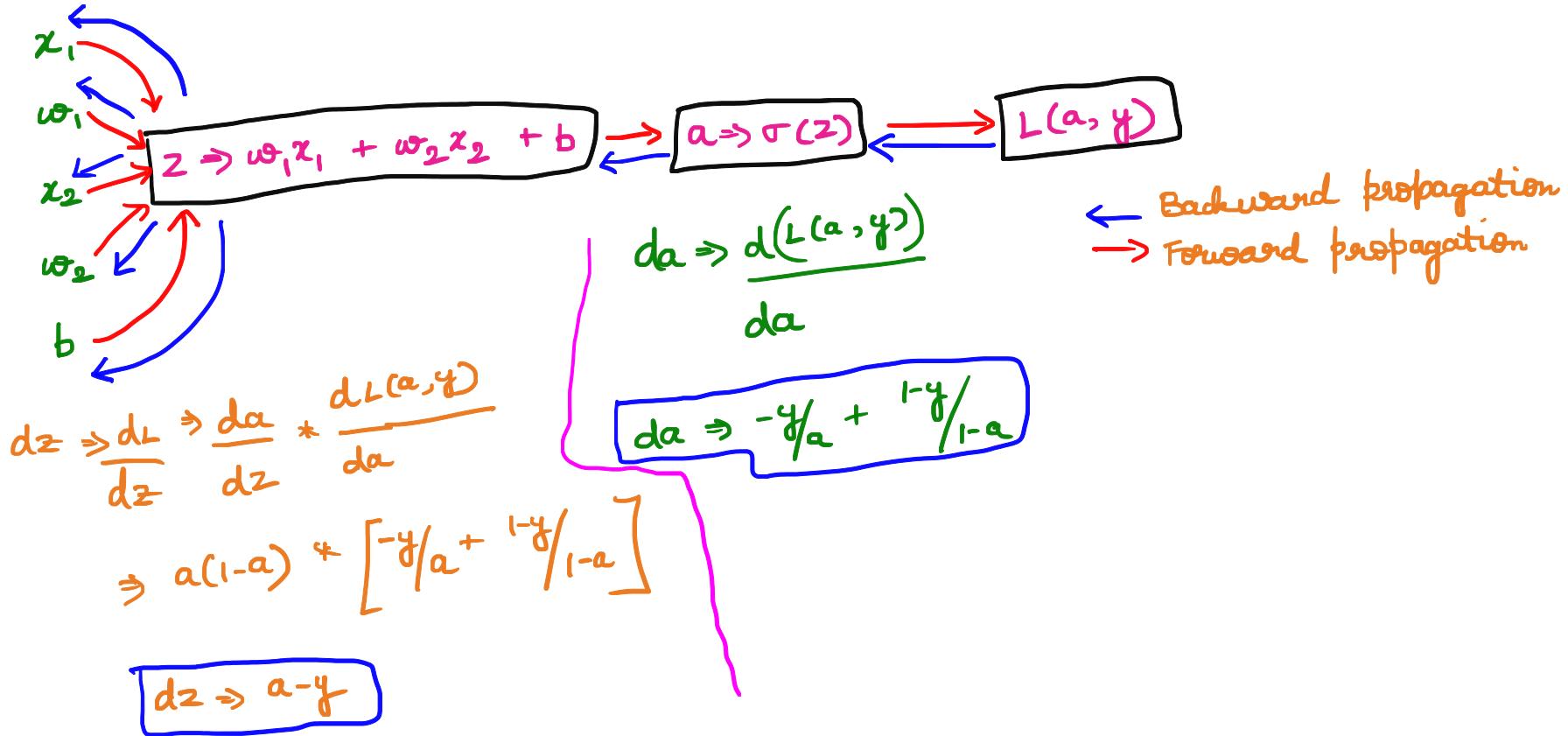
$$1) dz \Rightarrow \frac{dL}{dz} \Rightarrow a-y$$

$$\Rightarrow \frac{dL(a,y)}{da} \Rightarrow -y/a + 1-y/(1-a)$$

$$2) d\omega_1 \Rightarrow \frac{dL}{d\omega_1} \Rightarrow x_1 * dz$$

$$\Rightarrow \frac{da}{dz} \Rightarrow a + (1-a)$$

$$d\omega_2 \Rightarrow \frac{dL}{d\omega_2} \Rightarrow x_2 * dz$$



With these formulae, the slope and intercept are going to be updated in the gradient descent

$$\begin{aligned}
 w_1 &\Rightarrow w_1 - \alpha \frac{dL}{dw_1} \\
 w_2 &\Rightarrow w_2 - \alpha \frac{dL}{dw_2} \\
 b &\Rightarrow b - \alpha \frac{dL}{db}
 \end{aligned}$$

1) Derivative with respect to the Activation function

$$\frac{dL(a,y)}{da} \Rightarrow ?$$

$$L(a,y) \Rightarrow -\frac{1}{m} \left[\sum_{i=1}^m y \ln(a) + (1-y) \ln(1-a) \right]$$

$$\frac{d(L(a,y))}{da} \Rightarrow -y \ln(a) - (1-y) \ln(1-a)$$

$$\frac{da}{da} \Rightarrow -y \left| \frac{1}{a} \right. - (1-y) \left| \frac{1}{1-a} \right. + (-1)$$

$$\boxed{\frac{d(L(a,y))}{da} \Rightarrow -y/a + (1-y)/(1-a)}$$

2) Derivative with respect to sigmoid function

$$\frac{da}{dz} \rightarrow ?$$

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

$$a = \frac{1}{1+e^{-(w_1x_1 + w_2x_2 + b)}}$$

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

$$\begin{aligned}\frac{da}{dz} &\Rightarrow -1(1+e^{-z})^{-1-1}(0+e^{-z})(-1) \\ &\Rightarrow -1(1+e^{-z})^{-2}(-e^{-z})\end{aligned}$$

$$\Rightarrow \frac{e^{-z}}{(1+e^{-z})^2}$$

Add $+1$ and -1

$$\Rightarrow \frac{1}{(1+e^{-z})} * \frac{1+e^{-z}-1}{(1+e^{-z})}$$

$$\boxed{\frac{d}{dx}(x^n) \Rightarrow x^{n-1}}$$

$$\frac{da}{dz} \Rightarrow \frac{1}{1+e^{-z}} * \frac{e^{-z}-1}{1+e^{-z}}$$

$$\Rightarrow \text{sig}(z) * \left[1 - \frac{1}{1+e^{-z}} \right] \rightarrow \frac{1+e^{-z}+1}{1+e^{-z}}$$

$$\frac{da}{dz} \Rightarrow \text{sig}(z) * 1 - \text{sig}(z)$$

$$\boxed{\frac{da}{dz} \Rightarrow a * (1-a)}$$

3) Combining 1) and 2) together

$$dz \Rightarrow \frac{dL}{dz} \Rightarrow \frac{da}{dz} * \frac{dL(a,y)}{da}$$

$$\Rightarrow [a + (1-a)] * \left[-\frac{y}{a} + \frac{1-y}{1-a} \right]$$

$$\Rightarrow a * \cancel{(1-a)} * \frac{(1-a)(-y) + a(1-y)}{\cancel{a + (1-a)}}$$

$$dz \Rightarrow \frac{dL}{dz} \Rightarrow -y + \cancel{ay} + a - \cancel{af} \Rightarrow a - y \quad \boxed{dz = a - y}$$

Derivatives with respect to weights

$$z = \omega_1 x_1 + \omega_2 x_2 + b$$

$$\frac{dL(a,y)}{dw} \Rightarrow ?$$

$$z = w^T x$$

$$\boxed{\frac{dz}{dw} \Rightarrow x}$$

$$\frac{dL(a,y)}{dw} \Rightarrow \frac{dL}{d\omega} \Rightarrow \frac{dz}{dw} * \frac{da}{dz} * \frac{dL}{da} \Rightarrow \frac{dz}{dw} * \frac{dL}{dz}$$

T P

$$\Rightarrow x + \frac{dL}{dz}$$

$$\left[dz \Rightarrow \frac{dL}{dz} \Rightarrow a-y \right]$$

$$\frac{dL}{d\omega} \Rightarrow x + dz$$

$$\boxed{\frac{dL}{d\omega_1} \Rightarrow d\omega_1 \Rightarrow x_1 * dz}$$

(θ^1)

$$\boxed{\frac{dL}{d\omega_2} \Rightarrow d\omega_2 \Rightarrow x_2 * dz}$$

(θ^2)

$$\boxed{\frac{dL}{d\omega_1} \Rightarrow d\omega_1 \Rightarrow x_1 * (a-y)}$$

$$\boxed{\frac{dL}{d\omega_2} \Rightarrow d\omega_2 \Rightarrow x_2 * (a-y)}$$

Derivatives with respect to weights:

$$\frac{dL(a, y)}{dw} \rightarrow ?$$

$$\boxed{\frac{dz}{dL} \Rightarrow a-y}$$

$$\frac{dL(a, y)}{dw} \Rightarrow \frac{dz}{dw} * \cancel{\frac{da}{dz}} * \cancel{\frac{dL}{da}} \Rightarrow \frac{dz}{dw} * \frac{dL}{dz}$$

$$\frac{dL(a, y)}{dw} \Rightarrow \frac{dz}{dw} * (a-y) \rightarrow ①$$

Let's take the cost function

$$② \rightarrow L(a, y) \Rightarrow -\gamma_m \left[\sum_{i=1}^m \left(y \ln(a) + (1-y) \ln(1-a) \right) \right]$$

$$\ln(a) \rightarrow \ln\left(\frac{1}{1+e^{-x}}\right)$$

$$\rightarrow \ln\left[\left(1+e^{-x}\right)^{-1}\right]$$

$$\boxed{\ln(a) \Rightarrow -\ln\left(1+e^{-x}\right)}$$

$$\ln(1-a) \rightarrow \ln\left[1 - \frac{1}{1+e^{-x}}\right]$$

$$\Rightarrow \ln\left[\frac{1+e^{-x}-1}{1+e^{-x}}\right] \Rightarrow \ln\left[\frac{e^{-x}}{1+e^{-x}}\right] \Rightarrow \ln(e^{-x}) - \ln(1+e^{-x})$$

$$\boxed{\ln(1-a) \Rightarrow -y - \ln(1+e^{-x})}$$

$$\boxed{\ln\left(\frac{a}{b}\right) \Rightarrow \ln(a) - \ln(b)}$$

ln and e are inverse functions

$$L(a, y) \Rightarrow -\frac{1}{m} \left[\sum_{i=1}^m (y_i \ln(a) + (1-y_i) \ln(1-a)) \right] \quad \text{Note: We will use the outer -ve at the last}$$

$$\Rightarrow \sum_{i=1}^m \left[y_i [-\ln(1+e^{-x})] + (1-y_i) [-y_i - \ln(1+e^{-x})] \right]$$

$$\Rightarrow \sum_{i=1}^m \left[-y_i \ln(1+e^{-x}) - y_i - \ln(1+e^{-x}) + y_i x + y_i \ln(1+e^{-x}) \right]$$

$$\Rightarrow -y_i - \ln(1+e^{-x}) + y_i x$$

$$\Rightarrow y_i x - [y_i + \ln(1+e^{-x})]$$

Applying log and e on y

$$\boxed{\ln(ab) \Rightarrow \ln(a) + \ln(b)}$$

$$\Rightarrow y_i x - [\ln(e^x) + \ln(1+e^{-x})]$$

$$\Rightarrow y_i x - [\ln[(e^x) \cdot (1+e^{-x})]]$$

$$\Rightarrow y_i x - [\ln[e^x(1) + e^x e^{-x}]]$$

$$\Rightarrow y_i x - [\ln[1+e^x]]$$

$$L \Rightarrow y\hat{y} - \ln(1 + e^{\hat{y}})$$

$$z \Rightarrow w_0x + b$$

$$L \Rightarrow y(w_0x + b) - \ln(1 + e^{w_0x + b})$$

Taking the derivative with respect to w

$$\frac{dL}{dw} \Rightarrow \frac{dL}{dw} \left[y(w_0x + b) - \ln(1 + e^{w_0x + b}) \right]$$

$$\frac{dL}{dw} \Rightarrow yx - \left(\frac{e^{w_0x + b}}{1 + e^{w_0x + b}} \cdot x \right) \Rightarrow yx - \left(\frac{x \cdot e^{w_0x + b}}{1 + e^{w_0x + b}} \right)$$

$$\frac{dL}{dw} \Rightarrow yx - x \operatorname{sig}(z)$$

$$\Rightarrow x[y - \operatorname{sig}(z)]$$

(-ve) on both sides

$$\frac{dL}{dw} \Rightarrow -x[y - \operatorname{sig}(z)]$$

$$\frac{dL}{dw} \Rightarrow x[\operatorname{sig}(z) - y]$$

$$\boxed{\frac{dL}{dw} \Rightarrow x(a - y)}$$

$$\frac{e^{w_0x + b}}{1 + e^{w_0x + b}} \Rightarrow \frac{1}{\frac{1 + e^{w_0x + b}}{e^{w_0x + b}}} \Rightarrow \frac{1}{\frac{1}{e^{w_0x + b}} + \frac{e^{w_0x + b}}{e^{w_0x + b}}}$$

$$\Rightarrow \frac{1}{\frac{1}{e^{w_0x + b}} + 1}$$

$$\Rightarrow \frac{1}{e^{-(w_0x + b)}} + 1 \Rightarrow \frac{1}{1 + e^{-(w_0x + b)}}$$

$$\boxed{\frac{e^{w_0x + b}}{1 + e^{w_0x + b}} \Rightarrow \frac{1}{1 + e^{-(-w_0x + b)}} \Rightarrow \operatorname{sig}(z)}$$

Cross Entropy Loss function

$$L(a, y) \Rightarrow -\frac{1}{m} \left[\sum_{i=1}^m y \ln(a) + (1-y) \ln(1-a) \right]$$

Using the -ve which was left at the beginning.

sigmoid
function

(concluding the above findings)

$$\frac{d_L}{d\omega} \Rightarrow d\omega \Rightarrow x * (\alpha - y) \Rightarrow x + dz$$

Since we have x_1 and x_2

$$\frac{d_L}{d\omega_1} \Rightarrow d\omega_1 \Rightarrow x_1 * (\alpha - y)$$

$$\frac{d_L}{d\omega_2} \Rightarrow d\omega_2 \Rightarrow x_2 * (\alpha - y)$$

Derivative with respect to bias:

$$\frac{dL(a, y)}{db} \Rightarrow ?$$

$$dz \Rightarrow \frac{dL}{dz} \Rightarrow a - y$$

$$\frac{dL}{db} \Rightarrow \frac{dz}{db} + \cancel{\frac{da}{dz}} + \cancel{\frac{dL}{da}} \Rightarrow \frac{dz}{db} + \frac{dL}{dz}$$

Let's calculate $\frac{dz}{db} \Rightarrow ?$



deeplearning.ai

Basics of Neural Network Programming

Gradient descent
on m examples

Logistic regression on m examples

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

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

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

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

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

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

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

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

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

$$J / m \leftarrow$$

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

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

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

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

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

Vectorization

$$J=0 ; d\omega_1 \Rightarrow \frac{dL}{d\omega_1} \rightarrow 0 ; d\omega_2 \Rightarrow \frac{dL}{d\omega_2} \rightarrow 0 ; db \Rightarrow \frac{dL}{db} \rightarrow 0 ; dz \Rightarrow \frac{dL}{dz}$$

number of
train examples

for $i=1$ to m {

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

$$\alpha^{(i)} = \sigma(z^{(i)}) \quad (\because \sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}})$$

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

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

number of features

for loop {

$$d\omega_1 \Rightarrow d\omega_1 + x_1^{(i)} dz^{(i)}$$

$$d\omega_2 \Rightarrow d\omega_2 + x_2^{(i)} dz^{(i)}$$

$$db \Rightarrow db + dz^{(i)}$$

Vectorization

To avoid this

As the # of features,
the complexity too
increases

Here only there are
2 features. Think
of n features

$$J = \frac{J}{m}, d\omega_1 \Rightarrow \frac{d\omega_1}{m}, d\omega_2 \Rightarrow \frac{d\omega_2}{m}; db \Rightarrow \frac{db}{m}$$

} Taking the averages

$$\omega_1 \Rightarrow \omega_1 - \lambda d\omega_1$$

$$\omega_2 \Rightarrow \omega_2 - \lambda d\omega_2$$

$$b \Rightarrow b - \lambda db$$

This step will be done after finding the
derivatives of both the slopes and intercept



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 += w[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.

Matrix

$$u = Av$$

Non-Vectorized

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

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

1) for i ...

2) for j ...

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

Vectorized

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

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])
```

Non-vectorized version

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

X: vectorized version

some of the examples

Logistic regression derivatives

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

$$d\omega = np.zeros((n_x, 1))$$

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 $j=1 \dots n_x$
 $d\omega_j := \dots$

$$\boxed{\begin{aligned} dw_1 &+= x_1^{(i)} dz^{(i)} \\ dw_2 &+= x_2^{(i)} dz^{(i)} \\ db &+= dz^{(i)} \end{aligned}} \quad | \quad n_x = 2$$

$$d\omega += x^{(i)} dz^{(i)}$$

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

$$d\omega /= m.$$



deeplearning.ai

Basics of Neural Network Programming

Vectorizing Logistic Regression

Vectorizing Logistic Regression

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

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

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

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

$\xrightarrow{\text{---}}$

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

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

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

$$\dots \underline{\underline{\omega^T x^{(m)} + b}}$$

$$\rightarrow \overline{z} = \text{np.dot}(\underline{\underline{\omega^T}}, \underline{\underline{X}}) + \underline{\underline{b}} \in \mathbb{R}^{(1,1)}$$

"Broadcasting"

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

$$\text{cost function} \rightarrow z^{(1)} = w^T x^{(1)} + b \quad | \quad z^{(2)} = w^T x^{(2)} + b \quad | \quad z^{(3)} = w^T x^{(3)} + b$$

$$\text{activation function} \rightarrow a^{(1)} = \sigma(z^{(1)}) \quad | \quad a^{(2)} = \sigma(z^{(2)}) \quad | \quad a^{(3)} = \sigma(z^{(3)}) \quad | \quad \dots \quad M \text{ training examples}$$

1st training example

2nd training example

3rd training example

$$z \left[z^{(1)}, z^{(2)}, \dots, z^{(m)} \right] \xrightarrow{\text{X}} w^T \begin{bmatrix} \vdots & \vdots & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} + [b, b, \dots, b] \xrightarrow{(n_x, m) \hookrightarrow \mathbb{R}^{n \times m}}$$

$$\Rightarrow w^T X + [b, b, \dots, b]$$

$$\Rightarrow [w^T x^{(1)} + w^T x^{(2)} + \dots + w^T x^{(m)}] + [b, b, \dots, b] \xrightarrow{1 \times m}$$

$$z \left[z^{(1)}, z^{(2)}, \dots, z^{(m)} \right] \xrightarrow{z^{(1)} \quad z^{(2)} \quad \dots \quad z^{(m)}} \left[\underline{w^T x^{(1)} b} + \underline{w^T x^{(2)} b} + \dots + \underline{w^T x^{(m)} b} \right]$$

$$z \Rightarrow np.dot(w^T X) + b$$

Although b is a single value while adding with the vector multiplication, it spreads like the above. for addition

$$A = [a^{(1)}, a^{(2)}, \dots, a^{(m)}] = \sigma(z)$$

Python Broadcasting



deeplearning.ai

Basics of Neural Network Programming

Vectorizing Logistic Regression's Gradient Computation

Vectorizing Logistic Regression

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

$$dz^{(2)} = a^{(2)} - y^{(2)}$$

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

$1 \times m$

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

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

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

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

$$\begin{aligned} db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\ &= \frac{1}{m} \underbrace{\text{np. sum}(dZ)} \end{aligned}$$

$$\begin{aligned} dw &= \frac{1}{m} X dZ^T \\ &= \frac{1}{m} \left[\begin{array}{c|c} x^{(1)} & \dots & x^{(m)} \\ \hline 1 & & 1 \end{array} \right] \left[\begin{array}{c} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{array} \right] \\ &= \frac{1}{m} \left[\underbrace{x^{(1)} dz^{(1)}}{} + \dots + \underbrace{x^{(m)} dz^{(m)}}{} \right] \\ &\quad n \times 1 \end{aligned}$$

Activation function $\Rightarrow A \Rightarrow [a^{(1)}, a^{(2)}, \dots, a^{(m)}] \Rightarrow \sigma(z)$

Dependent variable $\Rightarrow Y \Rightarrow [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$

$dz \Rightarrow [dz^{(1)}, dz^{(2)}, dz^{(3)}, \dots, dz^{(m)}]$

$d\bar{z} \Rightarrow A - Y \Rightarrow [a^{(1)} - y^{(1)}; a^{(2)} - y^{(2)}; \dots, a^{(m)} - y^{(m)}]$

Looping over all the training examples

$$d\omega = 0$$

$$d\omega \rightarrow d\omega + x^{(1)} dz^{(1)}$$

$$d\omega \rightarrow d\omega + x^{(2)} dz^{(2)}$$

$$\vdots$$

$$d\omega \rightarrow d\omega/m$$

$$db = 0$$

$$db \rightarrow db + dz^{(1)}$$

$$db \rightarrow db + dz^{(2)}$$

$$\vdots$$

$$db \rightarrow db/m$$

Derivative of the intercept

Derivative of the slope

Vectorized format

Removed 1 for
Loop

We can exclude the for-loop by introducing vectorization

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

$$d\omega \rightarrow \frac{1}{m} X(dz)^T$$

$$\rightarrow \frac{1}{m} \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$d\omega \rightarrow \frac{1}{m} \left[x^{(1)} dz^{(1)} + x^{(2)} dz^{(2)} + \dots + x^{(m)} dz^{(m)} \right]$$

~~x~~: Vectorized Form ~~x~~

Implementing Logistic Regression

~~x~~: Non vectorized form ~~x~~

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

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 = sigmoid(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

From the data-set

L> 100 grams of Apple contains 56 calories from carbohydrates.

L> 100 grams of Beef contains 104 calories from proteins.

L> 100 grams of Eggs contains 99 calories from Fat.

52 calories from Protein.

4.4 calories from Carbohydrates.

Question

Calculate the % of calories from
carbohydrates, protein and fat for each food } ?

Let's take Apple

$$\begin{array}{r} \textcircled{1} \\ 56.0 \\ 1.2 \\ 1.8 \\ \hline 59.0 \end{array}$$

% of calories from carbohydrate $\Rightarrow \frac{56}{59} \Rightarrow 94.9\%$

protein $\Rightarrow \frac{1.2}{59.0} \Rightarrow 2\%$

fat $\Rightarrow \frac{1.8}{59.0} \Rightarrow 3\%$

Calculate all the
% of calories from
the foods

without using
explicit for loops

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

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$

$A = \begin{bmatrix} 56.0 & 0.0 & 4.4 & 68.0 \\ 1.2 & 104.0 & 52.0 & 8.0 \\ 1.8 & 135.0 & 99.0 & 0.9 \end{bmatrix}_{(3,4)}$

\downarrow^0

$\xrightarrow{1}$

59cal

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

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

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

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

$\uparrow^{(3,4)} / (1,4)$

Broadcasting example

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

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

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

Question - 1

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}_{m \times 1} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix}_{1 \times 4} \xrightarrow{\text{Expanding this 100}} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}_{m \times 1} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix}_{m \times 1} \Rightarrow \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

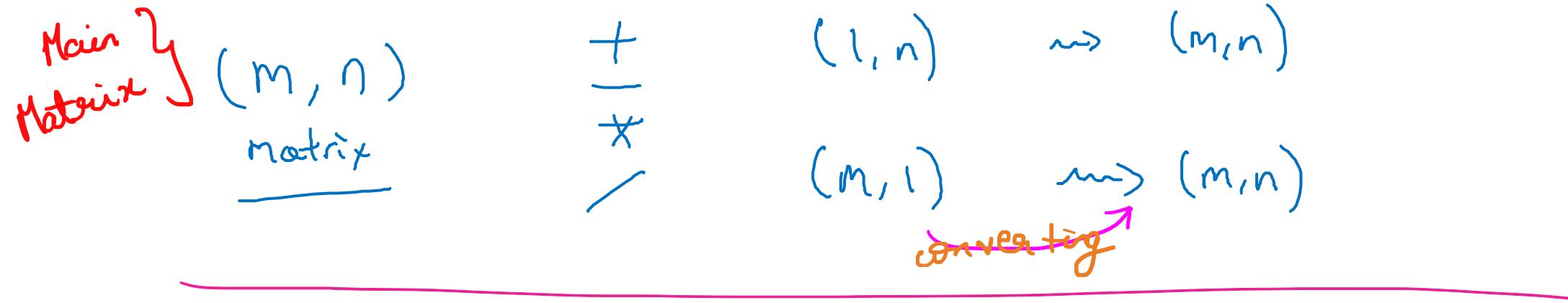
Question - 2

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{m \times n} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}_{1 \times n} \xrightarrow{\quad} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{m \times n} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}_{n \times m} \Rightarrow \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

Question - 3

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{m \times n} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}_{m \times 1} \xrightarrow{\quad} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{m \times n} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix}_{m \times n} \Rightarrow \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

General Principle



$$\begin{bmatrix} m, 1 \\ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \end{bmatrix} + \begin{matrix} R \\ 100 \end{matrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$
$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} + 100 = [101 \quad 102 \quad 103]$$

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

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

$p(y|x)$

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

←

$$\text{If } y=1: \quad p(y|x) = \hat{y} \underset{=} {=} 1$$
$$\text{If } y=0: \quad p(y|x) = \hat{y}^0 (1-\hat{y})^{(1-y)} = 1 \times (1-\hat{y}) = 1 - \hat{y}$$
$$\uparrow \log p(y|x) = \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y})$$
$$= -\frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)$$

↓

$$P(y|x) \rightarrow (\hat{y})^y (1-\hat{y})^{(1-y)}$$

$$\text{If } y \rightarrow 1 \rightarrow P(y|x) \rightarrow (\hat{y})^1 (1-\hat{y})^{(1-1)} \rightarrow \hat{y} (1-\hat{y})^0 \rightarrow \hat{y}$$

$$y \rightarrow 0 \rightarrow P(y|x) \rightarrow (\hat{y})^0 (1-\hat{y})^{(1-0)} \rightarrow 1 + (1-\hat{y}) \rightarrow 1-\hat{y}$$

positive example $\rightarrow P(y|x) \rightarrow \hat{y}$

Negative example $\rightarrow P(y|x) \rightarrow 1-\hat{y}$

proof

applying log on both sides

$$\log [P(y|x)] \Rightarrow \log [(\hat{y})^y (1-\hat{y})^{(1-y)}]$$

$$\log(ab) \rightarrow \log(a) + \log(b)$$

$$\Rightarrow \log(\hat{y}^y) + \log(1-\hat{y})^{(1-y)}$$

$$\log[P(y|x)] \Rightarrow y \log(\hat{y}) + (1-y) \log(1-\hat{y}) \Rightarrow -L(\hat{y}, y)$$

Logistic
Regression

While training an
algorithm we want
the probabilities
to be large

Minimizing the loss is equal to
Maximizing the probability.

$$\log[P(y|x)] \Rightarrow -\text{Loss}(\hat{y}, y)$$

Cost on m examples

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

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

Maximum likelihood estimation \nearrow

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

(minimize)

For minimizing the cost function, we are using the Maximum Likelihood Estimation