

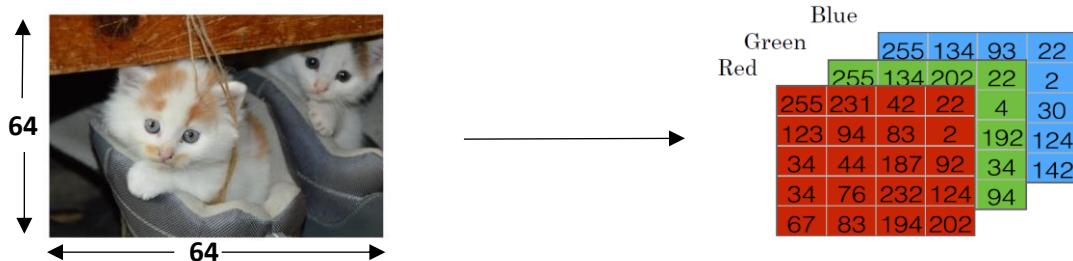
Binary Classification

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 $n_x = 64 \times 64 \times 3 = 12,288$.

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

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.

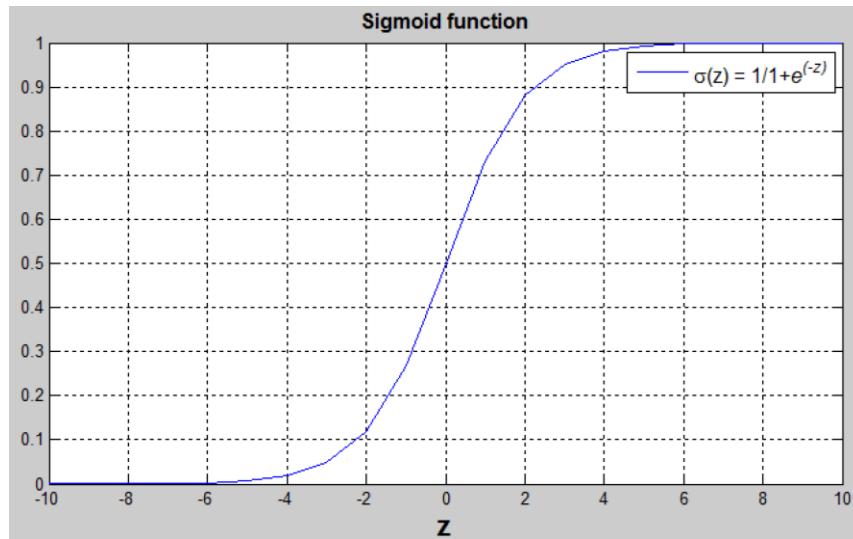
Example: Cat vs No - cat

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 $\sigma(z) = 0.5$

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

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.

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

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

- 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.

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



deeplearning.ai

Basics of Neural Network

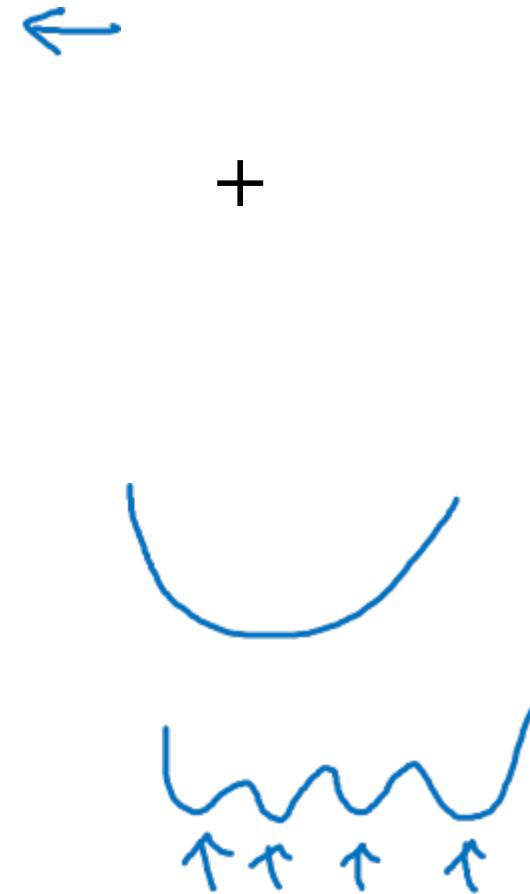
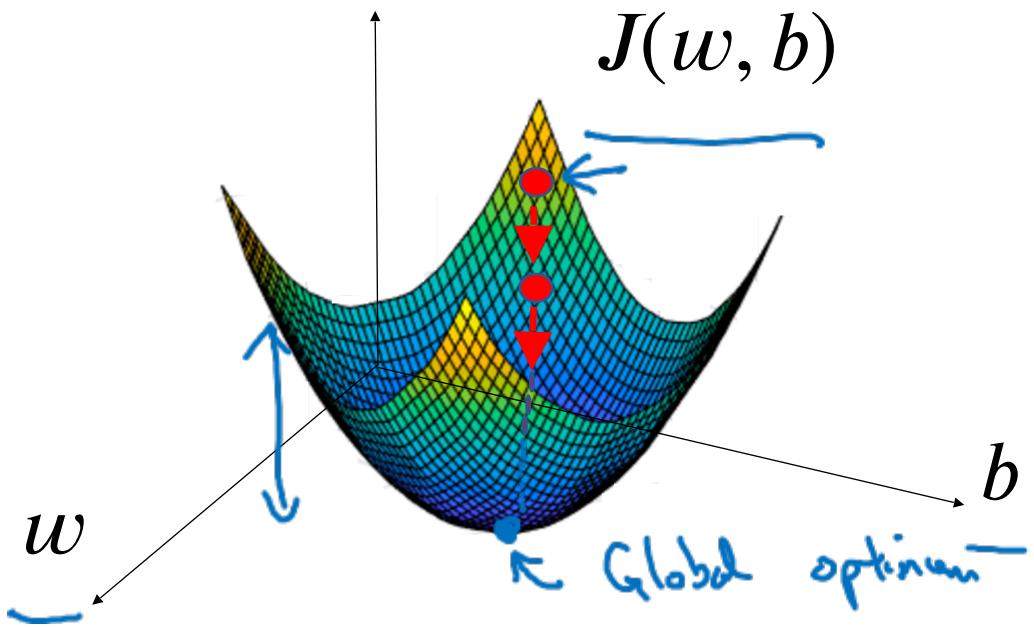
Programming Gradient Descent

Gradient Descent

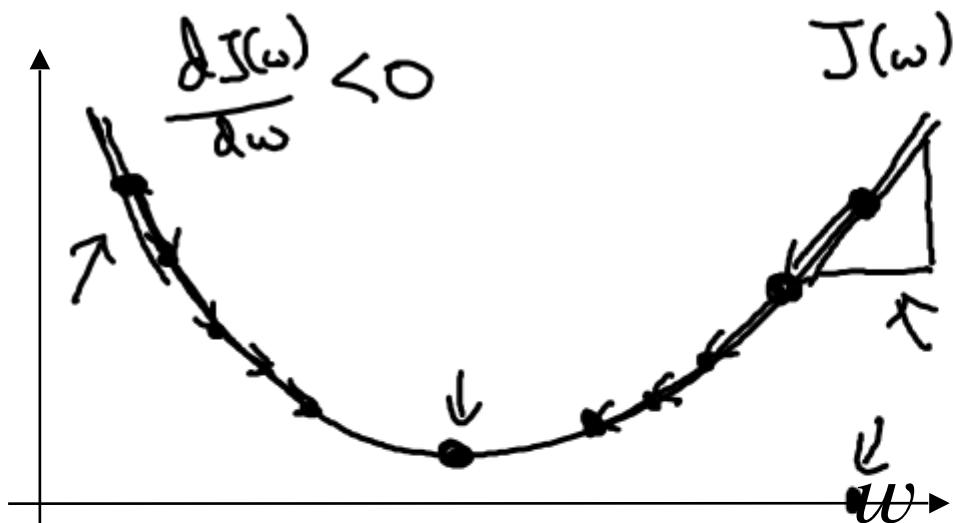
Recap: ,

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \text{[blue boxes]} - \frac{1}{m} \sum_{i=1}^m$$

Want to find that minimize



Gradient Descent



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

$$\underbrace{\frac{\partial J(w)}{\partial w}}_{=} ?$$

$$J(w, b)$$

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

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

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

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

"partial
derivative"
 J

$d w$
 $d b$

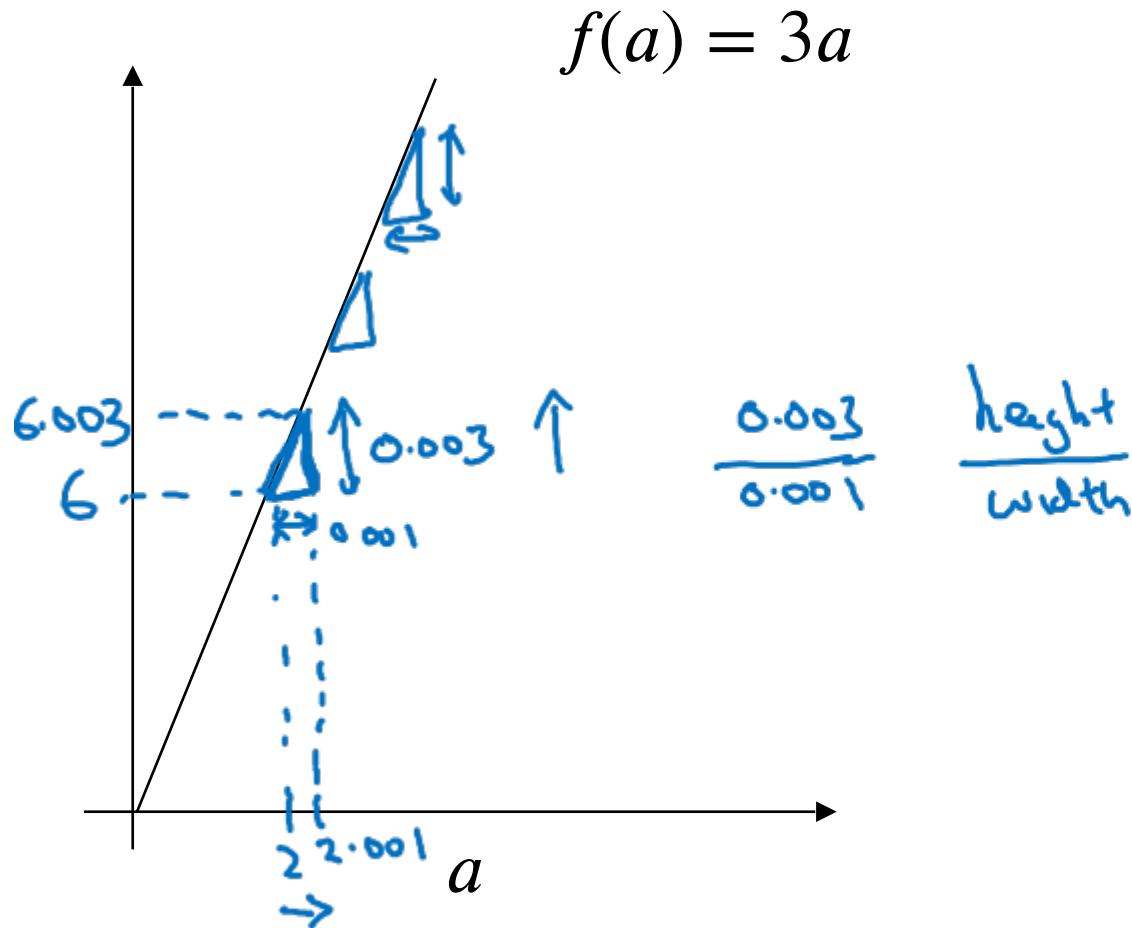


deeplearning.ai

Basics of Neural Network

Programming Derivatives

Intuition about derivatives



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

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

\curvearrowright slope (derivative) of $f(a)$
at $a=2$ is 3

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

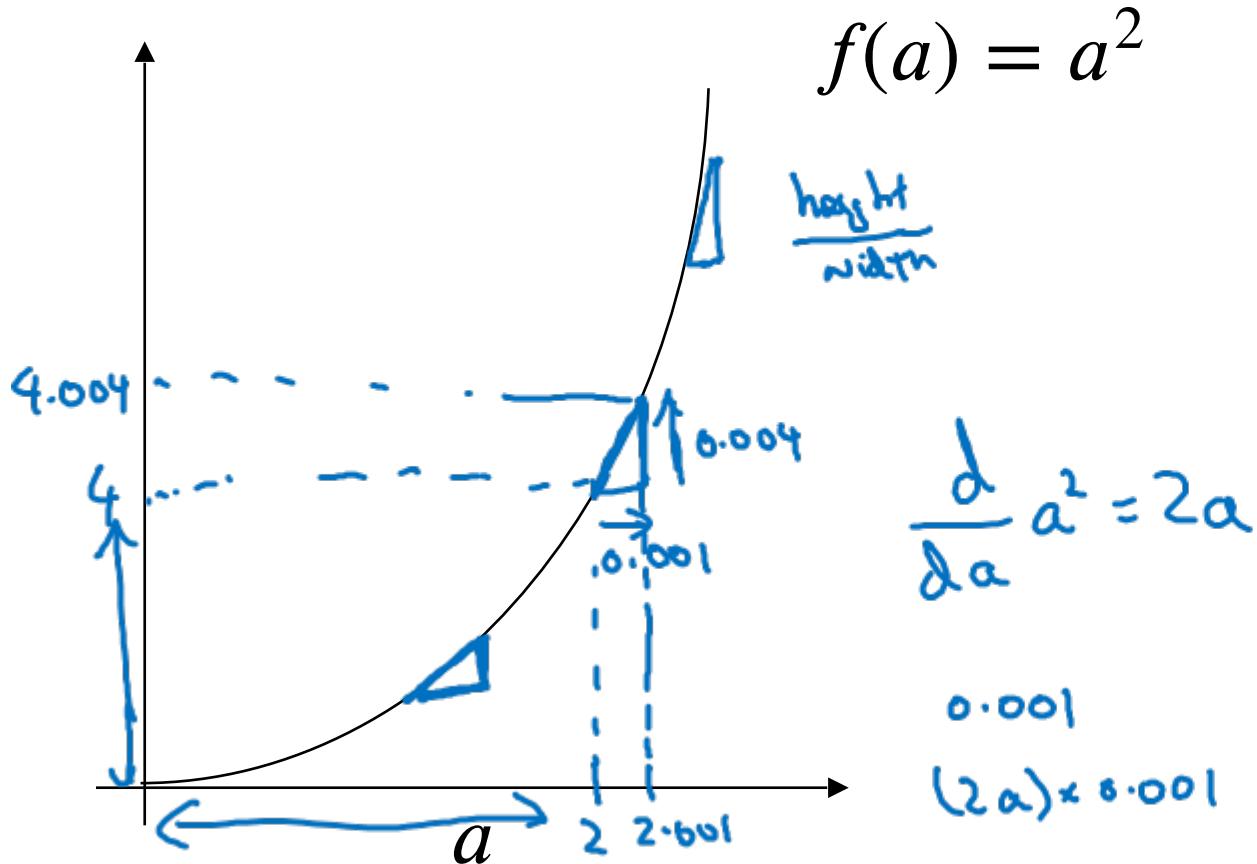
$a = 5.001 \quad f(a) = 15.003$

slope at $a=5$ is also 3

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

$0.001 \leftarrow$
 0.00000001
 0.0000000001

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.

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

$a = 5$

$a = 5.001$

$f(a) = 25$

$f(a) \approx 25.010$

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

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

More derivative examples

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

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

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

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

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

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

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

$$a = 2$$

$$a = 2.001$$

$$f(a) = 4$$

$$f(a) \approx 4.004$$

$$a = 2$$

$$a = \underline{2.001}$$

$$f(a) = 8$$

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

$$a = 2$$

$$a = \underline{2.001}$$

$$f(a) \approx 0.49315$$

$$f(a) \approx \underline{0.49265}$$

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



deeplearning.ai

Basics of Neural Network

Programming Computation Graph

Computation Graph

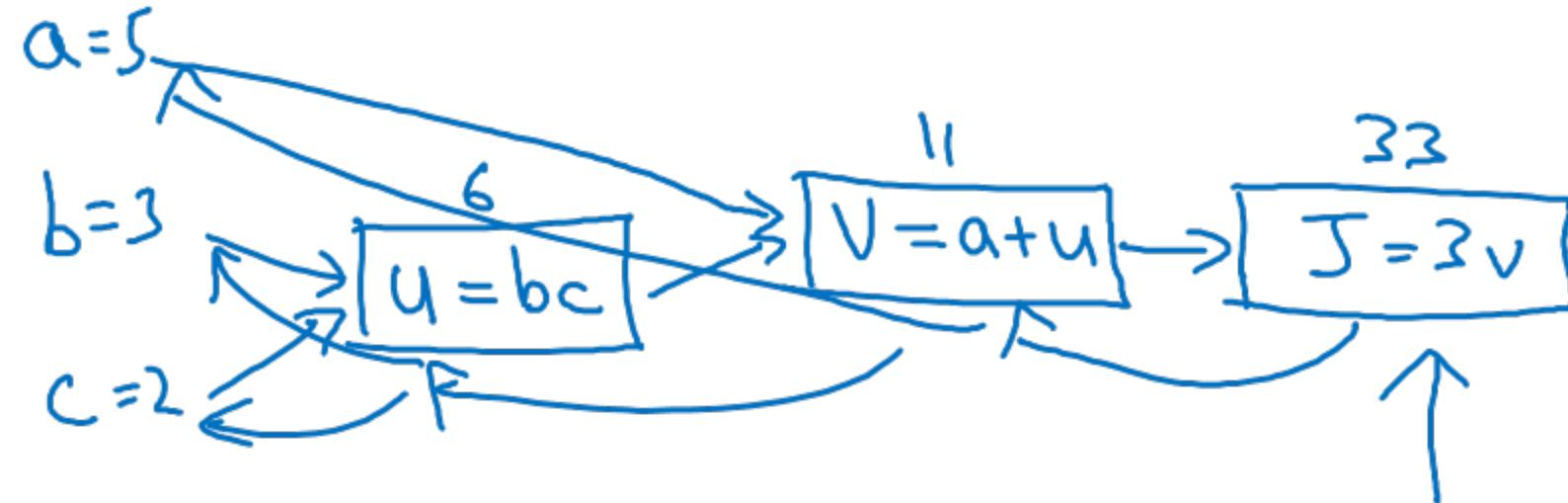
$$J(a, b, c) = 3(a + \underbrace{bc}_u) = 3(5 + 3 \times 2) = 33$$

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

$$u = bc$$

$$v = a + u$$

$$J = 3v$$





deeplearning.ai

Basics of Neural Network

Programming Derivatives with a Computation Graph

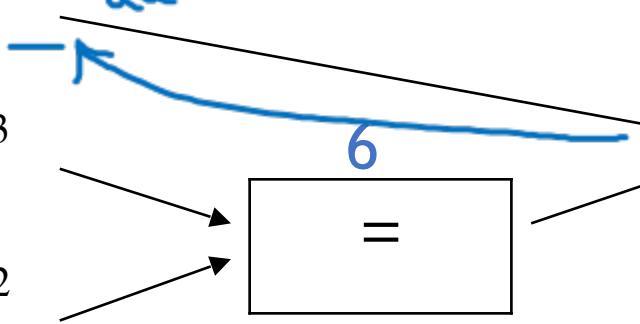
Computing derivatives

$$a = 5 \quad \frac{dJ}{da} \text{ "da" = 3}$$

$$b = 3$$

$$c = 2$$

$$=$$



$$11$$

$$33$$

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

$$J = 3v$$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$a \rightarrow v \rightarrow J$$

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

$$\frac{dJ}{da} = 3 = \frac{dJ}{dv} \frac{dv}{da} \quad 3 \times 1$$

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

$$\frac{\partial \text{FinalOutputVar}}{\partial \text{var}}$$

$$\frac{\partial J}{\partial \text{var}} \text{ "dvar"}$$

$$f(a) = 3a$$

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

$$J = 3v$$

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

Computing derivatives

$$a = 5$$

$$\frac{\partial J}{\partial a} \rightarrow \underline{\underline{\frac{\partial a}{\partial b} = 3}} \quad b = 3$$

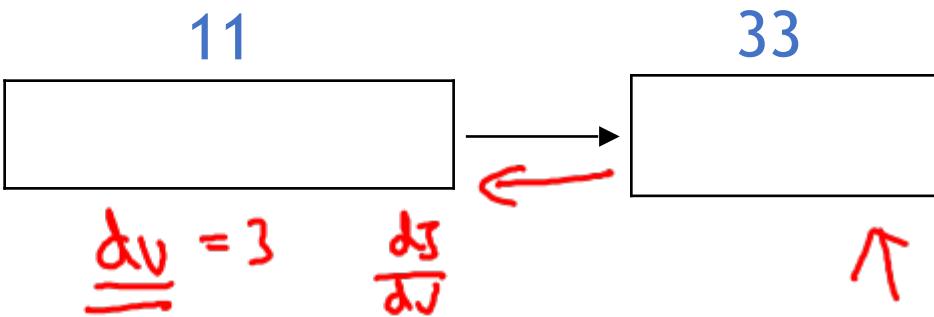
$$\frac{\partial J}{\partial b} \rightarrow \underline{\underline{\frac{\partial b}{\partial c} = 6}} \quad c = 2$$

$$\rightarrow \underline{\underline{\frac{\partial c}{\partial J} = 9}}$$

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

$$\underline{\underline{\frac{\partial J}{\partial b} = \frac{\frac{\partial J}{\partial u}}{\frac{\partial u}{\partial b}} = 6}} \quad \frac{\partial u}{\partial b} = 2$$

$$\underline{\underline{\frac{\partial J}{\partial a} = \left(\frac{\frac{\partial J}{\partial u}}{\frac{\partial u}{\partial a}} \right) \cdot \frac{\frac{\partial u}{\partial a}}{3} = 9}} \quad \frac{\partial u}{\partial a} = 3$$



$$u = 6 \rightarrow 6.001$$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$b = 3 \rightarrow 3.001$$

$$u = b \cdot c = 6 \rightarrow 6.002$$

$$c = 2 \\ 006$$

$$J = 33.006$$

$$V = 11.002 \\ J = 3V$$



deeplearning.ai

Basics of Neural
Network

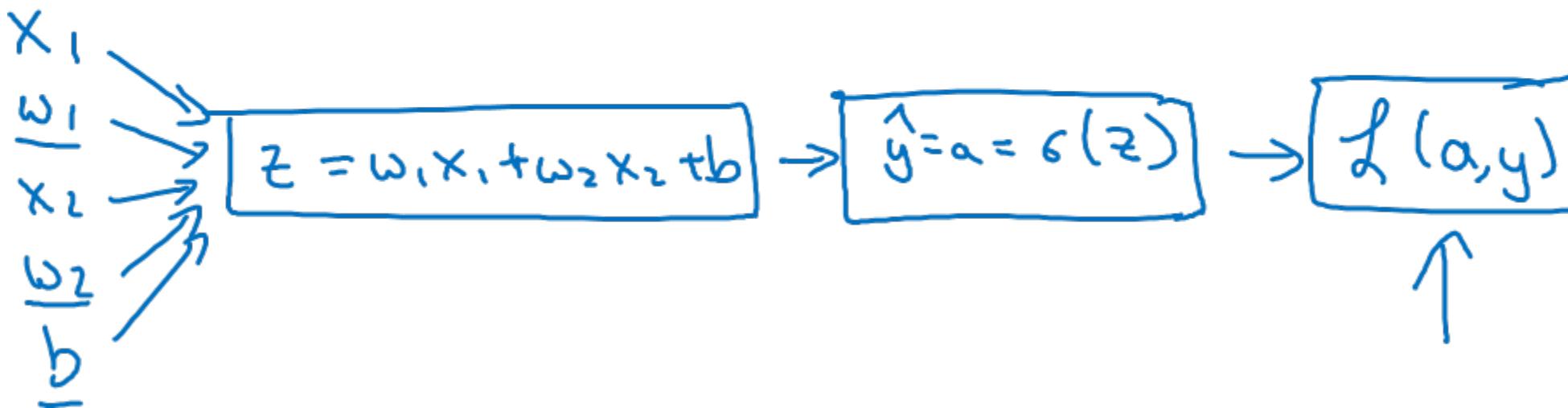
Programming
Logistic Regression
Gradient descent

Logistic regression recap

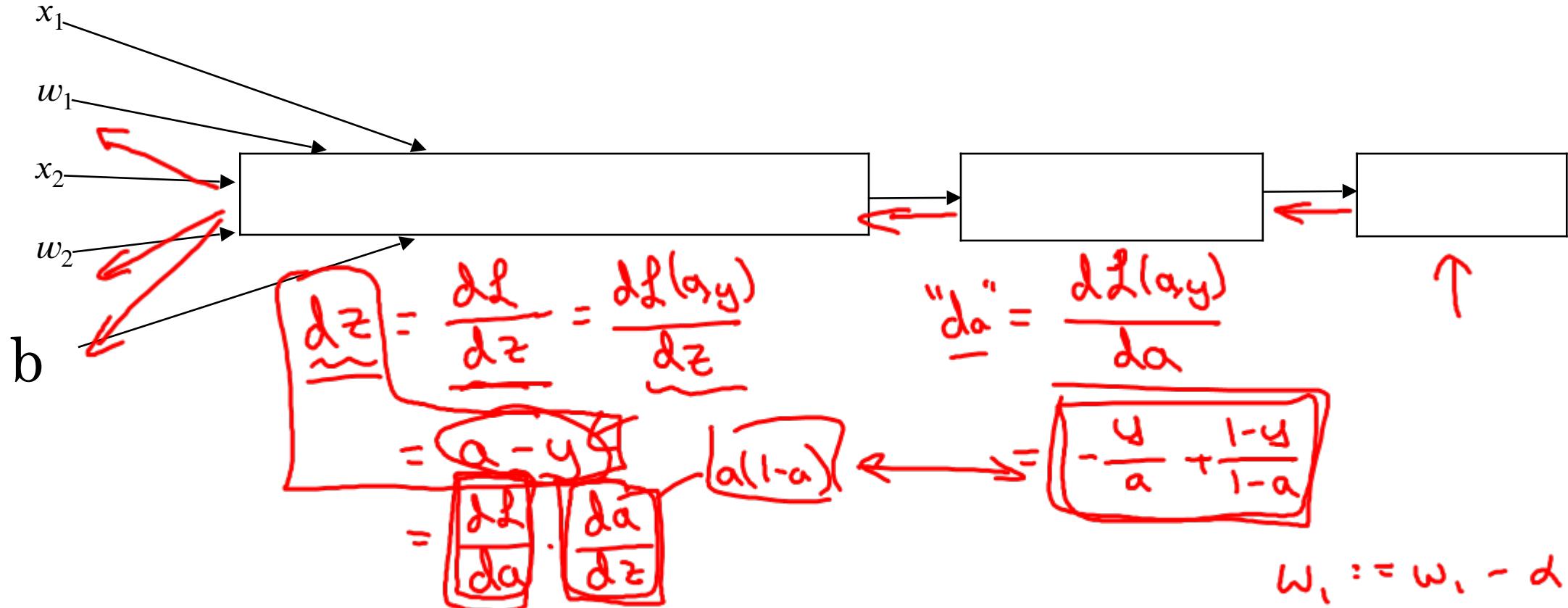
$$\rightarrow z = w^T x + b$$

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

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



Logistic regression derivatives



$$\frac{\partial L}{\partial w_1} = "d_{w_1}" = x_1 \cdot dz.$$

$$d_{w_1} = x_1 \cdot dz. \quad d_{b_1} = dz.$$

$$\begin{aligned} w_1 &:= w_1 - \alpha \frac{\partial L}{\partial w_1}, \\ w_2 &:= w_2 - \alpha \frac{\partial L}{\partial w_2}, \\ b &:= b - \alpha \frac{\partial L}{\partial b}. \end{aligned}$$



deeplearning.ai

Basics of Neural Network

Programming Gradient descent on m examples

Logistic regression on m examples

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

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

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

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

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

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

$$\underline{\frac{\partial z^{(i)}}{\partial w_i}} = \alpha^{(i)} - y^{(i)}$$

$$\left[\begin{array}{l} \frac{\partial J_t}{\partial w_1} = x_1^{(i)} \frac{\partial z^{(i)}}{\partial w_1} \\ \frac{\partial J_t}{\partial w_2} = x_2^{(i)} \frac{\partial z^{(i)}}{\partial w_2} \end{array} \right] \quad n=2$$

$$\left[\begin{array}{l} \frac{\partial J_t}{\partial b} = \frac{\partial z^{(i)}}{\partial b} \end{array} \right]$$

$$J / m \leftarrow$$

$$\frac{\partial J_t}{\partial w_1} / m; \frac{\partial J_t}{\partial w_2} / m; \frac{\partial J_t}{\partial b} / m. \leftarrow$$

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

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

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

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

Vectorization



deeplearning.ai

Basics of Neural Network

Programming Vectorization

What is vectorization?

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

Non-vectorized:

$$z = 0$$

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

$$z += b$$

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

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.

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

$J = 0, \boxed{dw1 = 0, dw2 = 0}, db = 0$

for i = 1 to n:

=

=

+=

d =

d +=

d +=

db += d

$n_x = 2$

$J = J/m, \boxed{d = d/m, d = d/m, db = db/m}$

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

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

$dw /= m.$



deeplearning.ai

Basics of Neural Network

Programming Vectorizing Logistic Regression

Vectorizing Logistic Regression

$$\begin{array}{l}
 \boxed{z^{(1)} = w^T x^{(1)} + b} \\
 \boxed{a^{(1)} = \sigma(z^{(1)})} \\
 \\
 \boxed{z^{(2)} = w^T x^{(2)} + b} \\
 \boxed{a^{(2)} = \sigma(z^{(2)})} \\
 \\
 \boxed{z^{(3)} = w^T x^{(3)} + b} \\
 \boxed{a^{(3)} = \sigma(z^{(3)})}
 \end{array}$$

$\overbrace{\mathbf{X}}^{\in R^{n_x \times n_m}} = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$
 $\overbrace{\mathbf{w}^T}^{\in R^{n_m \times 1}} = \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}$

$$\begin{aligned}
 \mathbf{z} &= \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \mathbf{w}^T \mathbf{X} + \begin{bmatrix} b & b & \dots & b \end{bmatrix}_{1 \times m} = \begin{bmatrix} w^T x^{(1)} + b \\ w^T x^{(2)} + b \\ \vdots \\ w^T x^{(m)} + b \end{bmatrix} \\
 &\rightarrow \mathbf{z} = \mathbf{np.dot}(\mathbf{w.T}, \mathbf{X}) + \underbrace{\mathbf{b}}_{(1, 1)} \quad \mathbf{R} \quad \text{"Broadcasting"}
 \end{aligned}$$

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



deeplearning.ai

Basics of Neural Network

Programming
Vectorizing Logistic
Regression's Gradient
Computation

Vectorizing Logistic Regression

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

$$\underline{dz} = [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}] \quad \leftarrow$$

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

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

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

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

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

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

$$dw = \frac{1}{m} \times dz'$$

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

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

Implementing Logistic Regression

$J = 0, d = 0, d = 0, db = 0$

$\boxed{\text{for } i = 1 \text{ to } m:}$

=

=

+=

$d =$

$[d +=$

$[d +=$

$db += d$

$J = J/m, d = d/m, d = d/m$

$db = db/m$

$$\left. \begin{array}{l} \\ \end{array} \right\} \delta w^t = X^{(i)} * \delta z^{(i)}$$

```
for iter in range(1000):  
    z = 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

Broadcasting example

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

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

$$= A_{(3,4)}$$



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

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

$$(m, n) \quad (2, 3) + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$$

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

$$(m, n) + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = \begin{bmatrix} 100 \\ 200 \end{bmatrix} \begin{bmatrix} 100 & 100 \\ 200 & 200 \end{bmatrix}$$

$(m, 1)$
 \downarrow
 (m, n)

\leftarrow
 \leftarrow

General Principle

$$\begin{array}{ccc} (m, n) & \xrightarrow{\quad \pm \quad} & (1, n) \rightsquigarrow (m, n) \\ \underline{\text{matrix}} & \cancel{*} & (m, 1) \rightsquigarrow (m, n) \end{array}$$

$$\begin{array}{ccccc} (m, 1) & + & R & & \\ \left[\begin{smallmatrix} 1 \\ 2 \\ 3 \end{smallmatrix} \right] & + & 100 & = & \left[\begin{smallmatrix} 101 \\ 102 \\ 103 \end{smallmatrix} \right] \\ [1 \ 2 \ 3] & + & 100 & = & [101 \ 102 \ 103] \end{array}$$

Matlab/Octave: bsxfun



deeplearning.ai

Basics of Neural Network

Programming A note on python/ numpy vectors

Python Demo

Python / numpy vectors

```
import numpy as np  
  
a = np.random.randn(5)  
  
a = np.random.randn( (5,1) )  
  
a = np.random.randn( (1,5) )  
  
assert(a.shape = (5,1))
```