

This week is about basics such as loss & cost functions, gradient descent optimization, computation graph and its derivatives, vectorization (representation of the IN/OUT relation in vector form so to train the parameter on batches of data instead of iterating over the samples).



deeplearning.ai

Basics of Neural Network Programming

Binary Classification

These concepts are presented with a logistic regression example.

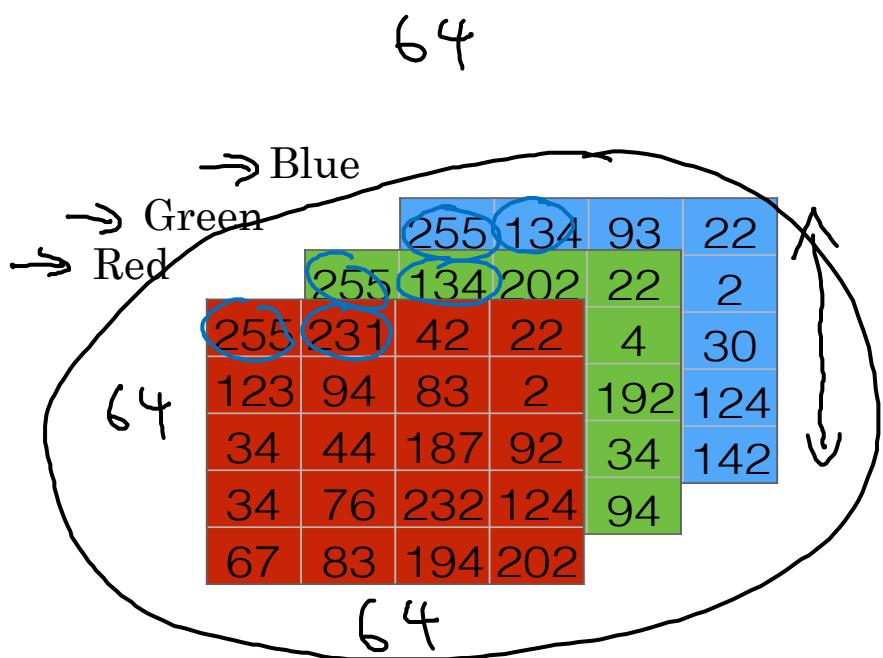
Binary Classification

= is there a cat in the picture?



1 (cat) vs 0 (non cat)

y



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

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

$$n = n_x = 12288$$

$$X \longrightarrow y$$

Notation

In NN notation convention, samples are stacked by column.

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

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

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

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

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

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

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

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

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

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



deeplearning.ai

Basics of Neural Network Programming

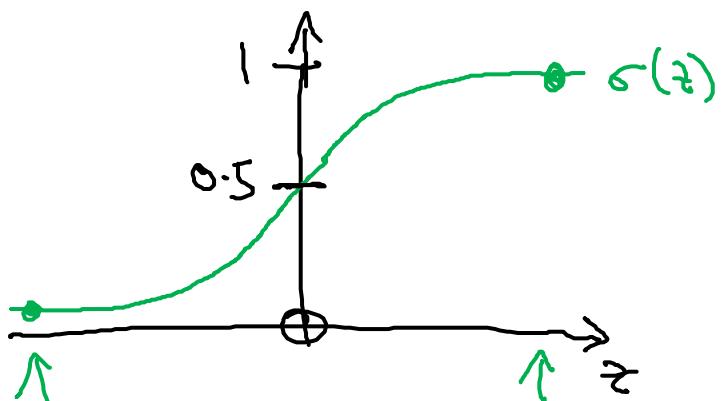
Logistic Regression

Logistic Regression

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

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

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



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

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

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_{n_x} \end{bmatrix} \quad \left. \begin{array}{l} \{w_0\} \rightarrow b \\ \{w_1, w_2, \dots, w_{n_x}\} \rightarrow w \end{array} \right.$$

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

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

If z large negative number

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



deeplearning.ai

Basics of Neural Network Programming

Logistic Regression cost function

Logistic Regression cost function

- cost function = average loss fn. over the entire data samples

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

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

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

i-th example.

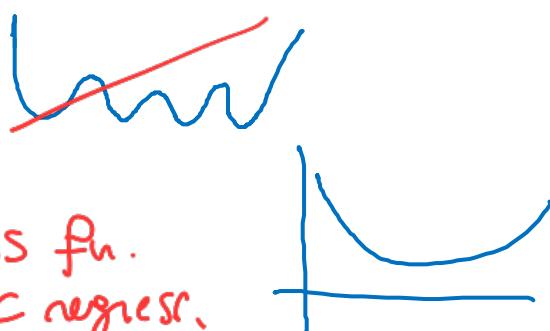
Loss (error) function:

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

term valid if $y=1$

term valid if $y=0$

This is the loss fn. used in logistic regress.



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



deeplearning.ai

Basics of Neural Network Programming

Gradient Descent

intuition, applied to logistic reg.

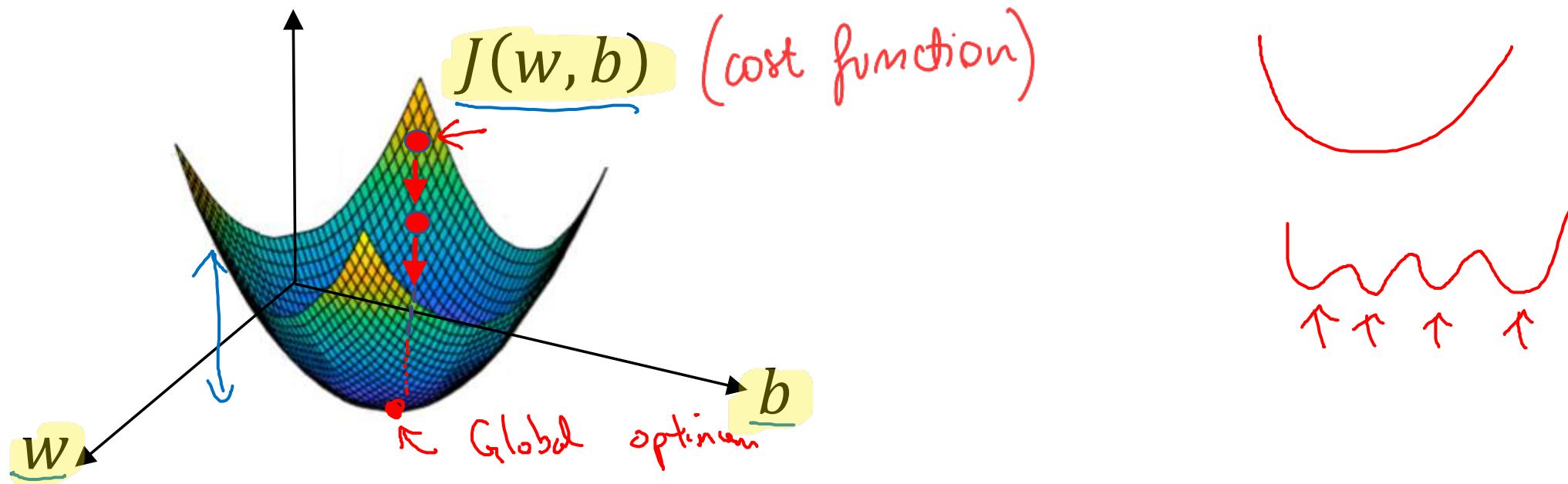
Gradient Descent

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

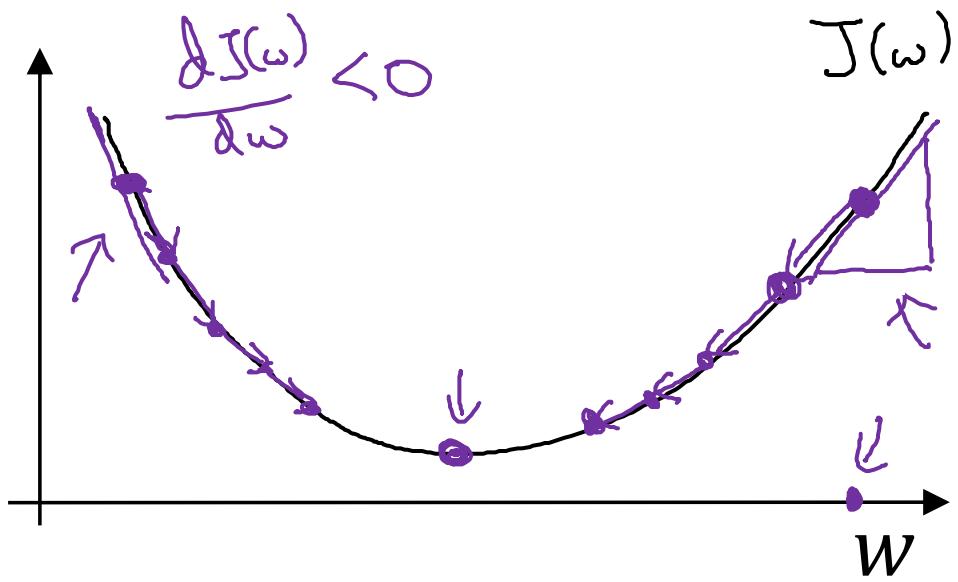
COST FUNCTION

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

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



Gradient Descent



$$J(w, b)$$

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

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

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

update rules of
gradient descent

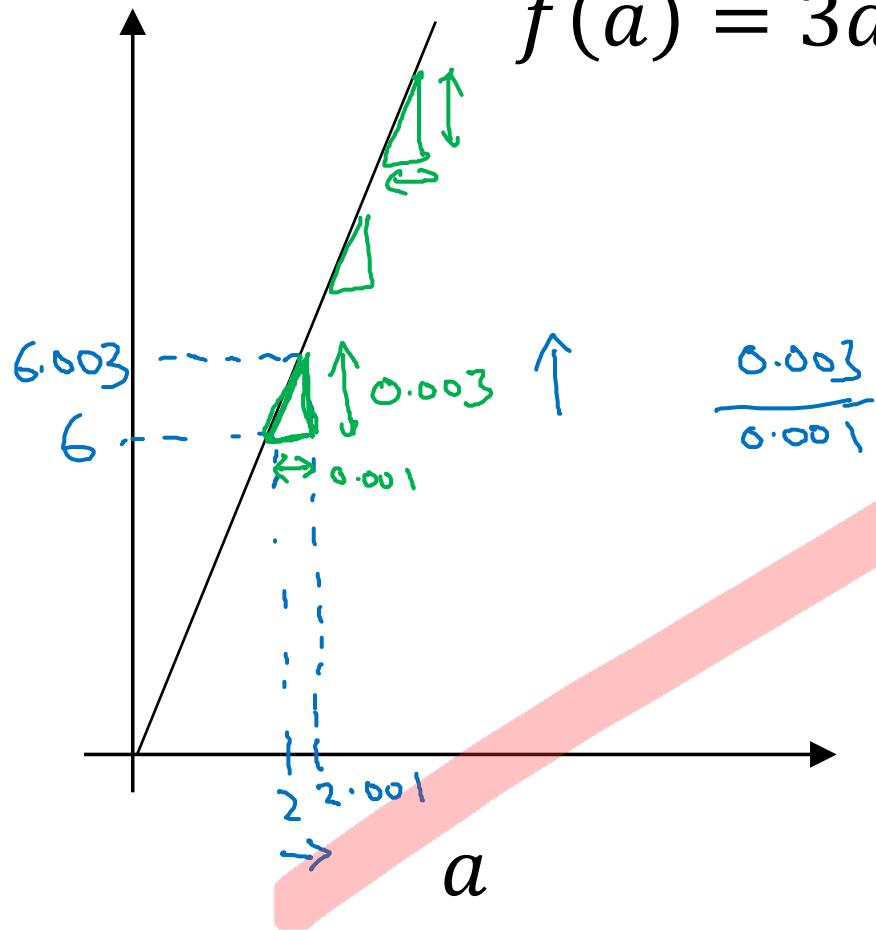


deeplearning.ai

Basics of Neural Network Programming

Derivatives

Intuition about derivatives



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

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

slope (derivative) of $f(a)$

at $a=2$ is 3

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

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

slope at $a=5$ is also 3

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

0.001 ←
0.00000001
0.0000000001

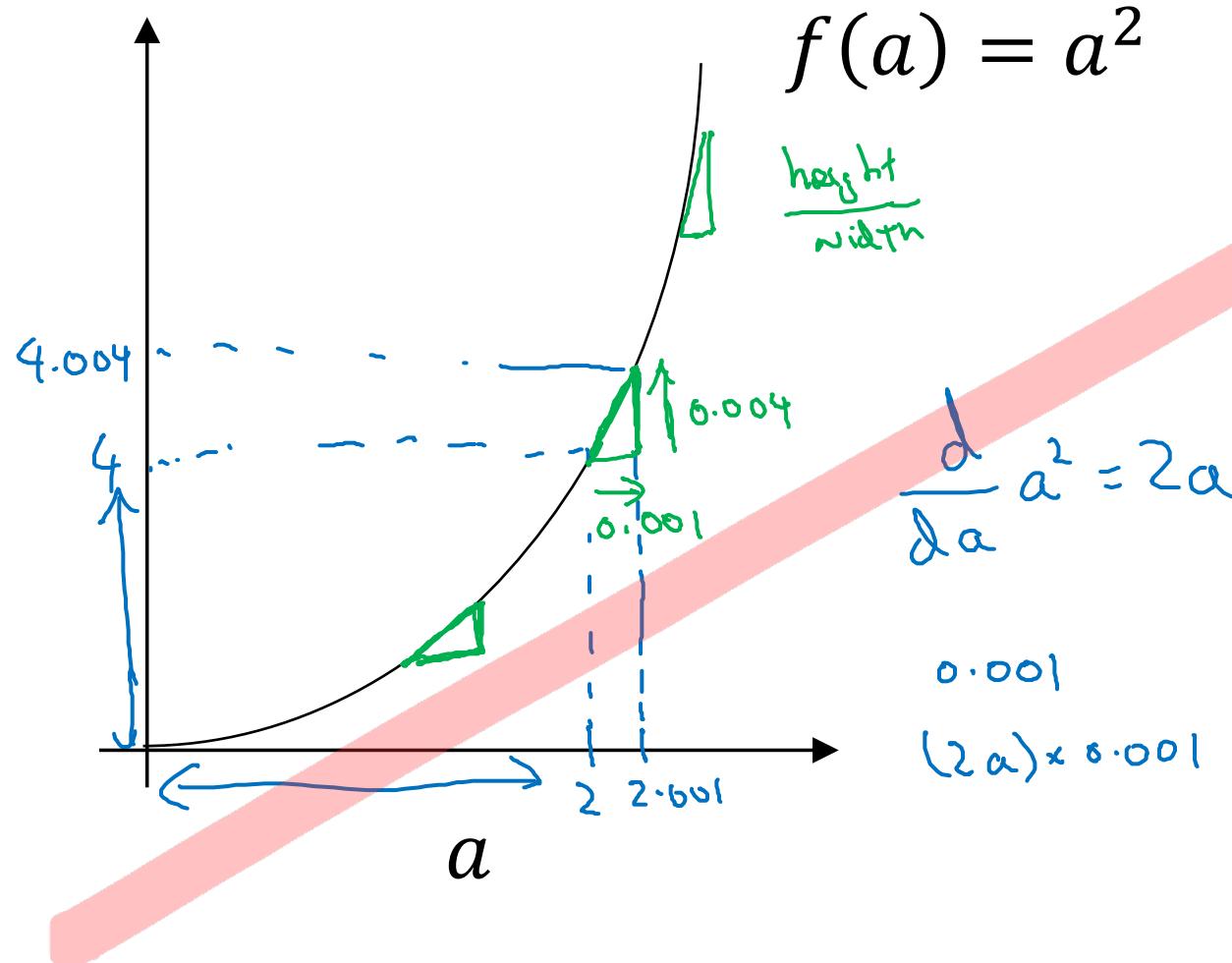


deeplearning.ai

Basics of Neural Network Programming

More derivatives
examples

Intuition about derivatives



$$a = 2$$

$$a = 2.001$$

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 $\boxed{a=5}$

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

More derivative examples

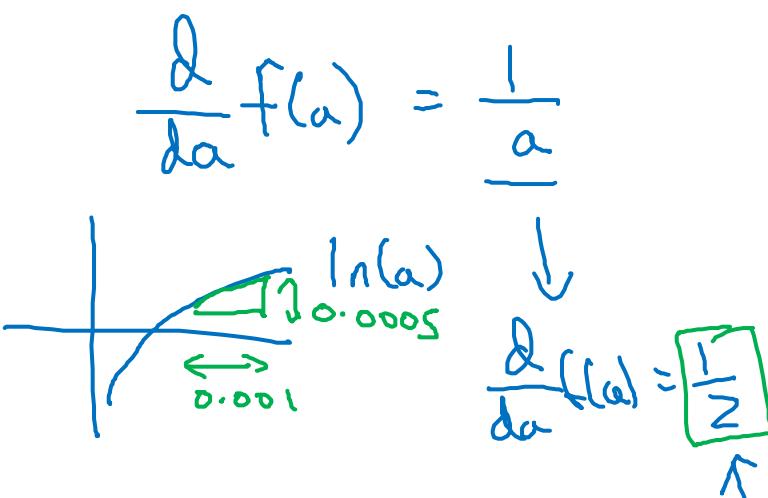
$$f(a) = a^2$$

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

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

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

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



$$a = 2$$

$$a = 2.001$$

$$f(a) = 4$$

$$f(a) \approx 4.004$$

$$a = 2$$

$$a = \underline{2.001}$$

$$f(a) = 8$$

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

$$a = 2$$

$$a = \underline{2.001}$$

$$f(a) \approx 0.69315$$

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

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



deeplearning.ai

Basics of Neural Network Programming

Computation Graph

- Computational graphs are structures used to represent a complex mathematical function as a combination of basic operations. Specifically, CGs are directed graphs in which the nodes are basic mathematical operations.
- Neural networks can be represented as CGs.

Computation Graph

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

$$\begin{array}{c} u \\ \downarrow \\ v \\ \downarrow \\ J \end{array}$$

$$u = bc$$

$$v = a + u$$

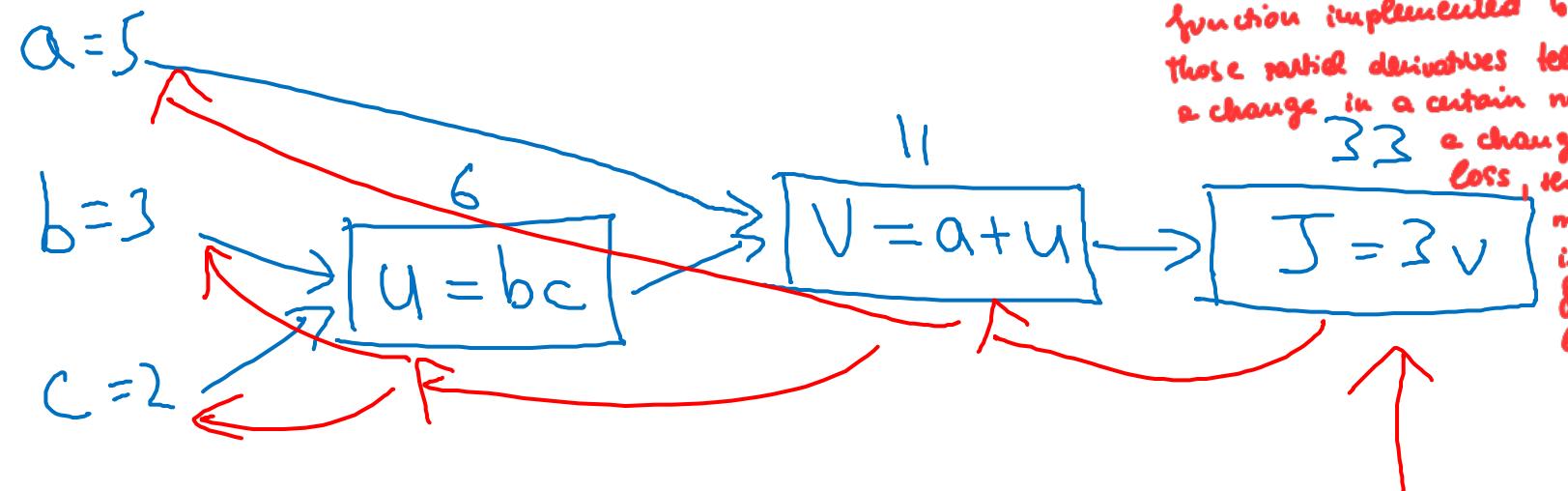
$$J = 3v$$

- Computations in a neural network are organized in a forward and a backward pass, or propagation.

FORWARD PROPAGATION: to compute the loss value $L(x, y)$, given an input $\{x_i\}$

BACKWARD PROPAGATION: to compute the gradient of the loss and distribute it across the nodes (neurons) that generated it. That is, to compute the gradient of the loss wrt the function implemented by each neuron.

Those partial derivatives tell how much a change in a certain neuron causes a change in the final loss, that is, how much each neuron is responsible for the overall loss.





deeplearning.ai

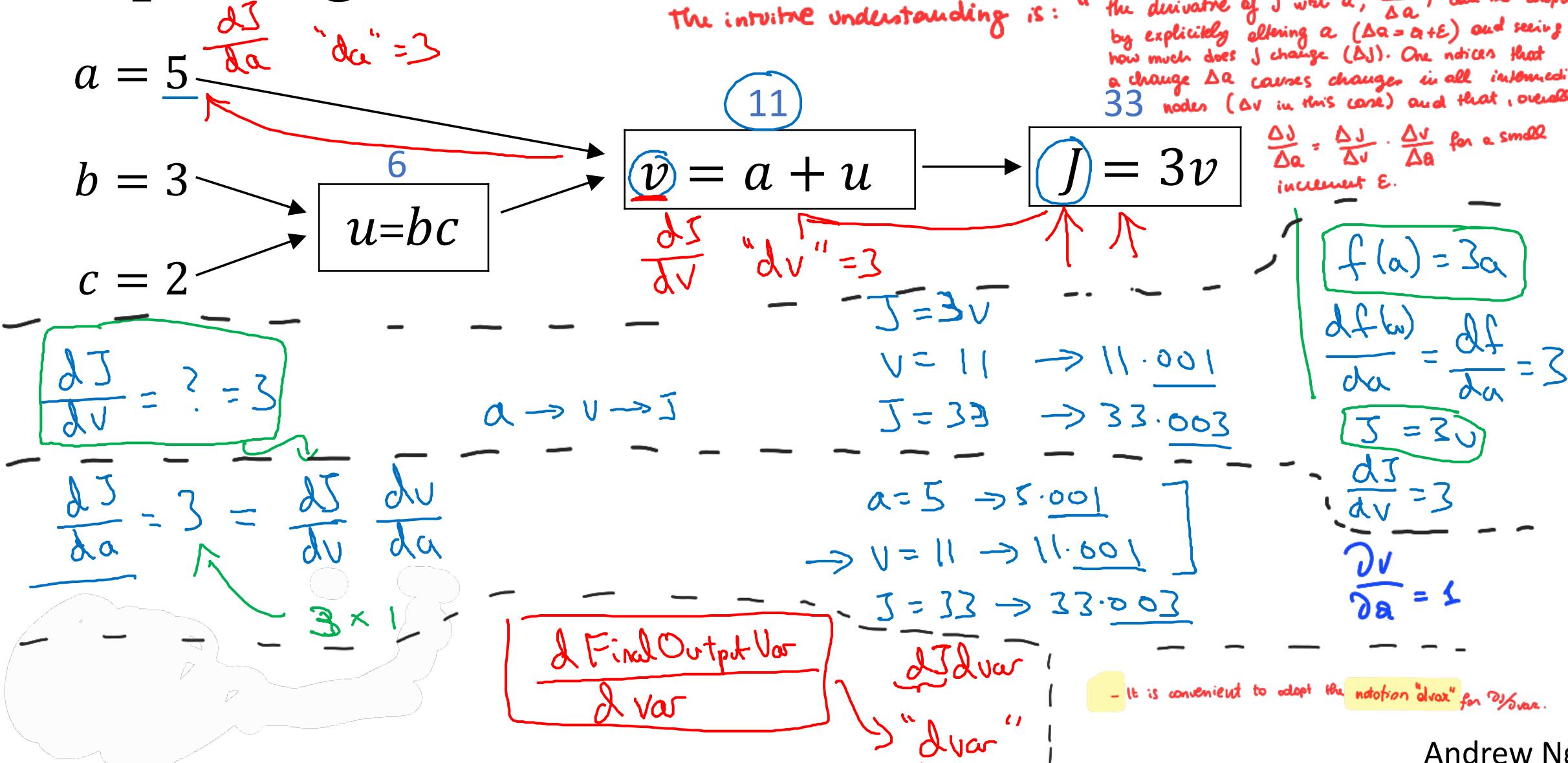
Basics of Neural Network Programming

Derivatives with a Computation Graph

Computing derivatives

- The computation graph gives an intuitive understanding of the **CHAIN RULE**: $\frac{\partial J(v(a))}{\partial a} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial a}$

The intuitive understanding is: "the derivative of J wrt a , $\frac{\Delta J}{\Delta a}$, can be computed by explicitly altering a ($\Delta a = a + E$) and seeing how much does J change (ΔJ). One notices that a change Δa causes changes in all intermediate nodes (Δv in this case) and that, overall,



Computing derivatives

$$\begin{aligned}
 a &= 5 \\
 \frac{\partial J}{\partial a} &\rightarrow \underline{\frac{\partial a}{\partial a}} = 3 \\
 b &= 3 \\
 \frac{\partial J}{\partial b} &\rightarrow \underline{\frac{\partial b}{\partial b}} = 6 \\
 c &= 2 \\
 \rightarrow \underline{\frac{\partial c}{\partial c}} &= 9 \\
 \frac{\partial J}{\partial u} &= 3 = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} \\
 &\quad \left(\begin{array}{c} 3 \\ -1 \end{array} \right)
 \end{aligned}$$

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

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

$$\begin{aligned}
 &11 \\
 &v = a + u \\
 &\underline{\frac{\partial v}{\partial v}} = 3 \quad \underline{\frac{\partial J}{\partial J}} \\
 &J = 3v
 \end{aligned}$$

$$\begin{aligned}
 u &= 6 \rightarrow 6.001 \\
 v &= 11 \rightarrow 11.001 \\
 J &= 33 \rightarrow 33.003
 \end{aligned}$$

$$\begin{aligned}
 b &= 3 \rightarrow 3.001 \\
 u &= b \cdot c = 6 \rightarrow 6.002 \\
 J &= 33.006
 \end{aligned}$$

$$\begin{aligned}
 v &= 11.002 \\
 J &= 3V
 \end{aligned}$$

- Since in NN training we need to compute $\frac{\partial J}{\partial v_{\text{out}}}$ for each var (=node) in the network, it is efficient to do so backwards (to reuse past computations). That is, it is convenient to compute $\frac{\partial J}{\partial v}$ before $\frac{\partial J}{\partial a} = \frac{\partial J}{\partial v} \frac{\partial v}{\partial a}$.

→ the BACKPROPAGATION algorithm refers to computing the partial derivative of the cost function wrt to (the function implemented by) each neuron starting from the last neurons and moving backward toward the first neurons.

→ It is efficient because it exploits the chain rule to reuse the partial derivatives of the later neurons in the computation of the derivatives of the former ones.

→ After computing the partial derivatives wrt all the nodes, the free parameters can be updated, e.g. via gradient descent. (in the following example (logistic regression) the free parameters are w_1, w_2, b).



deeplearning.ai

Basics of Neural Network Programming

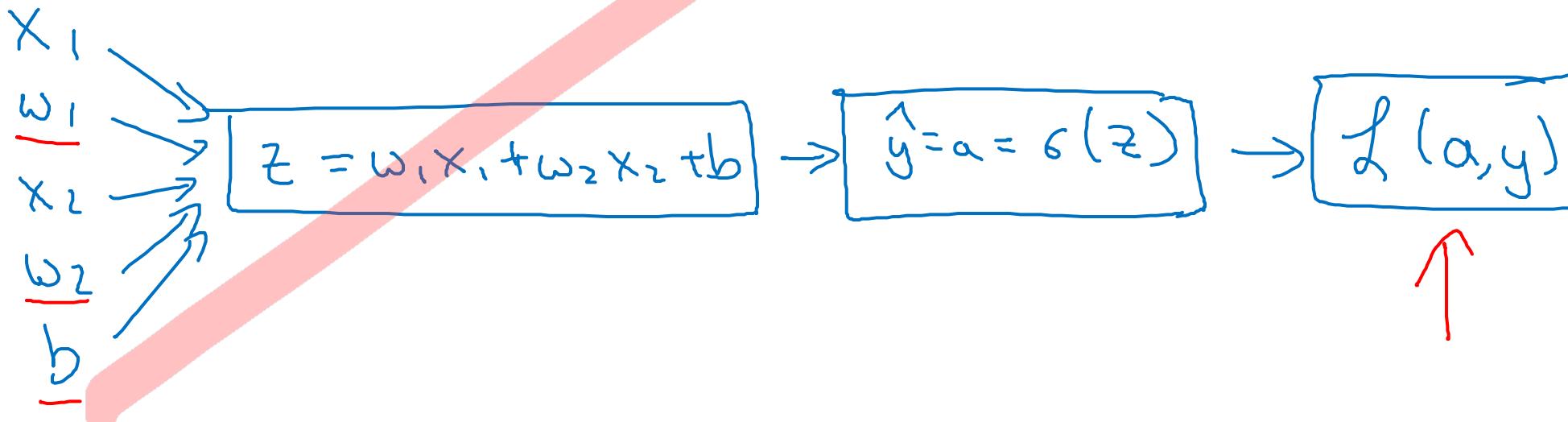
Logistic Regression Gradient descent

Logistic regression recap

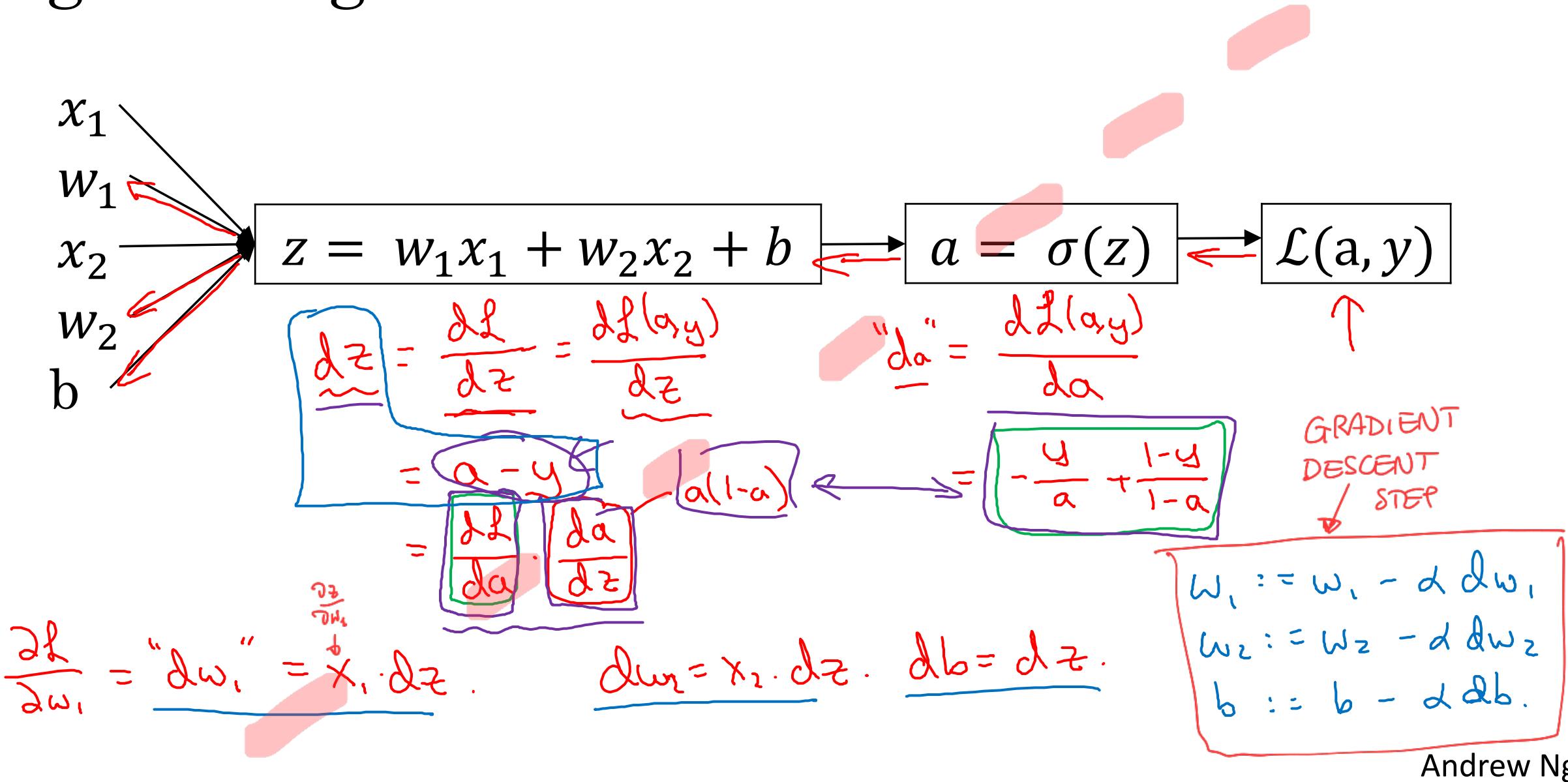
$$\rightarrow z = w^T x + b$$

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

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



Logistic regression derivatives



VECTORIZATION



deeplearning.ai

Basics of Neural Network Programming

Gradient descent
on m examples

Logistic regression on m examples

If we have more than one sample

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m l(a^{(i)}, y^{(i)})$$

$$\rightarrow a^{(i)} = \hat{y}^{(i)} = g(z^{(i)}) = g(\omega^\top x^{(i)} + b)$$

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

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

$$\frac{\partial}{\partial w_j} J(\omega, b) =$$

$$\frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_j} l(a^{(i)}, y^{(i)})$$

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

Being the cost function the
average loss over all samples
(and being the average a weighted sum)
then also the gradient of the cost
function wrt a neuron is the average of
the gradient of the cost function wrt that
neuron over all samples.

IN PRACTICE, if we have many samples,
to apply gradient descent to one variable
we need to:

- compute the gradient of the loss wrt that variable for each sample $\rightarrow dw^{(i)}$
- average those gradients (ie. compute the gradient of the cost function wrt that variable) $\rightarrow dw = \frac{1}{m} \sum_{i=1}^m dw^{(i)}$
- update the variable $var \leftarrow var - \alpha \cdot dw$

Logistic regression on m examples

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

\rightarrow For $i = 1$ to m

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

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

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

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

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

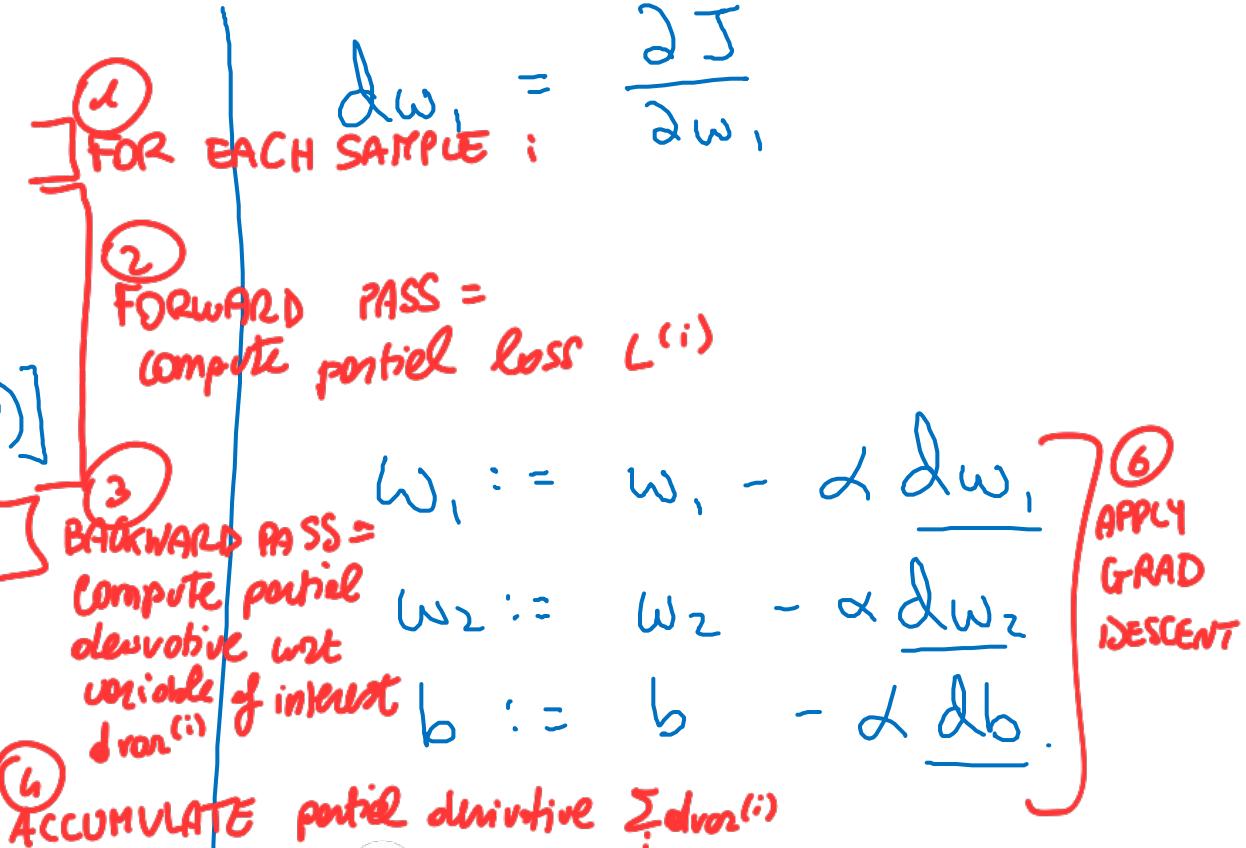
$$dw_2 += x_2^{(i)} dz^{(i)}$$

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

$$J / m \leftarrow$$

$$dw_1 / m; dw_2 / m; db / m.$$

⑤ Compute average $d\omega^{(i)}$ $\rightarrow \frac{\sum_{i=1}^m d\omega^{(i)}}{m}$





deeplearning.ai

Basics of Neural Network Programming

Vectorization

What is vectorization?

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

Non-vectorized:

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

$z += b$

SLOWER

- Replacing for loops with vector operations.
- Those operations are generally implemented in Python using GPU or CPU instructions that make use of parallelism, called SIMD instructions.

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

Vectorized

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

FASTER

\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.

$$u = Av$$

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

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

for i ...

 for j ...

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

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

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  
v/v
```

Logistic regression derivatives

How to vectorize the algorithm saw before for gradient descent for logistic regression.

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

~~for i = 1 to n:~~

$$z^{(i)} = w^T x^{(i)} + b \rightarrow z = \text{dot}(w^T, x) + b$$

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

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

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

for $j=1 \dots n$
 $dw_j += \frac{\partial J}{\partial w_j}$

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

$$n_x = 2$$

$$dw / m$$

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

$$dw = np.zeros((n_y \times [y_1, y_2, \dots, y_m]))$$

- hstack variables $w = [w_1, w_2], b = [b_1, b_2]$ $z = [z_1, z_2, \dots, z_m]$
- hstack outputs & gradients $dW = [dw_1, dw_2, \dots, dw_m]$
- hstack inputs $X = [x_1, x_2, \dots, x_m]$

~~$$dw += X^{(i)} dz^{(i)} \rightarrow dW = \text{dot}(X, dz^T)$$~~



deeplearning.ai

Basics of Neural Network Programming

Vectorizing Logistic Regression

Vectorizing Logistic Regression

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

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

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

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

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

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

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

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

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

"Broadcasting"

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



deeplearning.ai

Basics of Neural Network Programming

Vectorizing Logistic Regression's Gradient Computation

Vectorizing Logistic Regression

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

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

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

$1 \times m$

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

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

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

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

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

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

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 = 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

NOTE: In numpy do not use `rank=1` arrays (`shape = (n,)`)
because they do not behave consistently as row or column vectors.

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

59 cal

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

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

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

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

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

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

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

General Principle

$$\begin{array}{c} (m, n) \\ \text{matrix} \\ \hline + \\ \times \\ / \end{array} \quad \begin{array}{c} (1, n) \\ (m, 1) \end{array} \quad \rightsquigarrow (m, n)$$

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

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

EXPLAINED WILLY NILLY, you can IGNORE

Logistic regression cost function

- which is equivalent to writing

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

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

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

For some reason one might want to maximizing thus

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

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

and since \log is monotonically increasing, the above maxim. corresponds to maximizing thus

$$\begin{aligned} \log p(y|x) &= \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ &= -\frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) \end{aligned}$$

- The prediction of logistic regression is defined as:
 $\hat{y} = P(y=1|x)$

$$p(y|x)$$

These can be summarized as

Cost on m examples

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

$$\log p(\text{---}) = \sum_{i=1}^m \underbrace{\log p(y^{(i)} | x^{(i)})}_{\text{Maximum likelihood estimation}} \quad \nearrow$$

$$- L(\hat{y}^{(i)}, y^{(i)})$$

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

Cost: $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \quad \leftarrow \text{OR MINIMIZING THE NEGATIVE LIKELIHOOD}$