

①

[60 #15]

→ compute derivative

→ let see computational graph for logistic regression algorithm.

& do the forward & backward pass on the graph.

→ using graph we can compute gradient / derivative / slope efficiently.

Logistic Regression Recap:

$$① \quad Z = \omega^T x + b$$

$$② \quad \hat{y} = a = \sigma(Z)$$

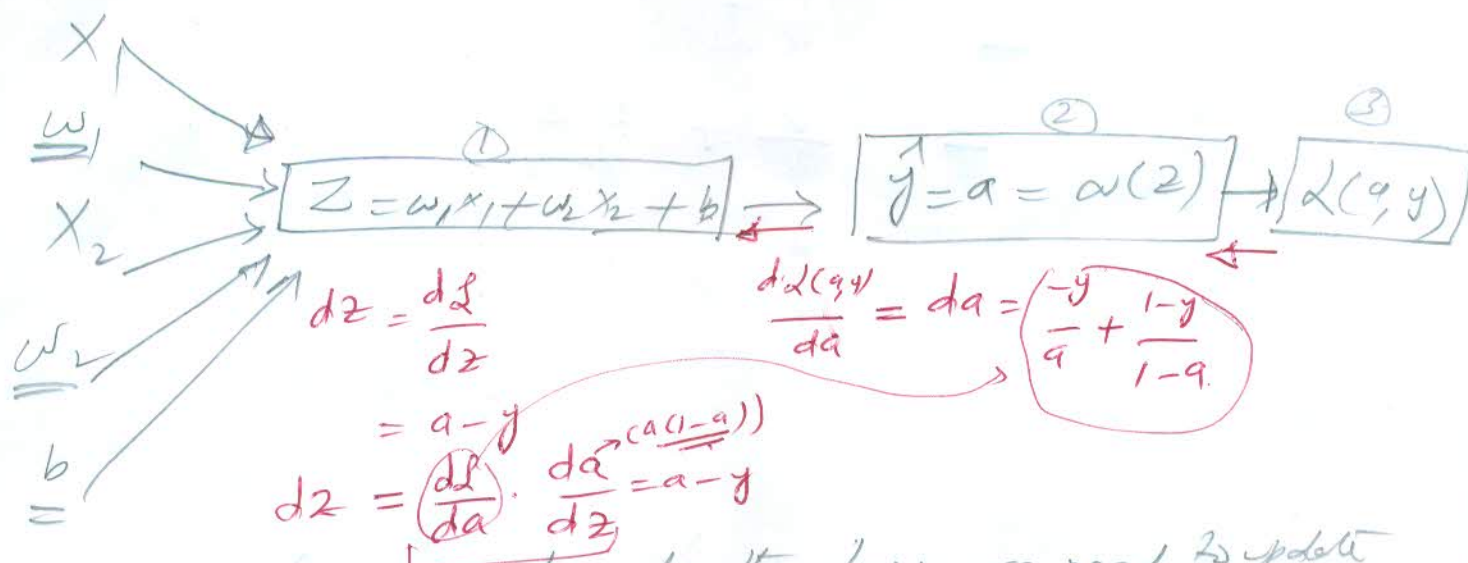
$$③ \quad \mathcal{L}(a, y) = -[y \log(a) + (1-y) \log(1-a)]$$

(for 1 example)

$$\frac{\omega^T x}{[---] [1]}$$

$$Z = \theta^T x + \theta_0$$

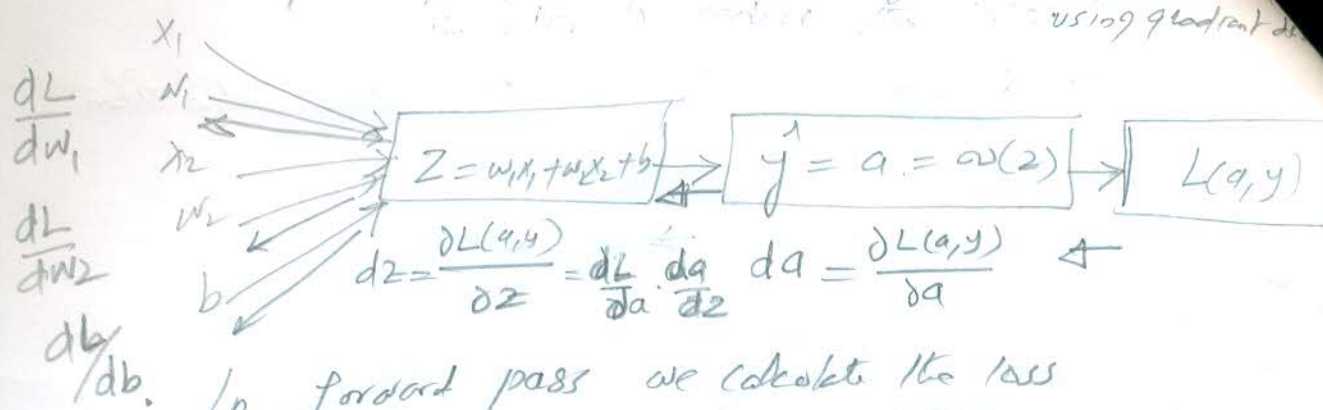
$$h_\theta(x) = \hat{y} = a = \sigma(Z)$$



In order to reduce the loss, we need to update ω_1, ω_2 & b — so we need to compute the derivative

2

- In order to Reduce the loss ($L(a, y)$) for a single Training example, we have to modify the parameters using gradient descent.



In forward pass we calculate the loss

In Backward pass we calculate the gradients

Use those gradients & ^{update} modify the parameters (w_1, w_2, b) to reduce the loss

$$\Rightarrow \frac{\partial L(a, y)}{\partial a} = ?$$

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

$$= -\frac{y}{a} + (1-y) \cdot \frac{1}{1-a} \cdot (-1)$$

$$\frac{\partial L(a, y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a}$$

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

$$\Rightarrow \frac{\partial L(a, y)}{\partial Z} = ?$$

$$\frac{\partial L(a, y)}{\partial Z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial Z} = da \cdot \left(\frac{\partial a}{\partial Z} \right)$$

(Chain Rule)

As $a = \frac{1}{1 + e^{-Z}}$

$$\boxed{\frac{\partial a}{\partial Z} = a(1-a)}$$

// property of sigmoid function derivative

③

$$dz = \frac{\partial L}{\partial z} = \left(-\frac{y}{a} + \frac{1-y}{1-a} \right) \times a(1-a)$$

$$= -y(1-a) + (1-y)a$$

$$= -y + ya + a - ya$$

$$\boxed{dz = a - y}$$

$$dw_1 = \frac{\partial L(a, y)}{\partial w_1} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_1} = dz \cdot \left(\frac{\partial z}{\partial w_1} \right)$$

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

$$\boxed{\frac{\partial z}{\partial w_1} = x_1}$$

$$\therefore \boxed{dw_1 = (a - y) \cdot x_1}$$

$$dw_2 = \frac{\partial L(a, y)}{\partial w_2} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_2} = dz \cdot \frac{\partial z}{\partial w_2}$$

$$a) \quad z = w_1 x_1 + w_2 x_2 + b$$

$$\boxed{\frac{\partial z}{\partial w_2} = x_2}$$

$$\boxed{dw_2 = (a - y) x_2}$$

$$= \frac{1}{a} \frac{dz}{db} = 1$$

(4)
$$db = \frac{\partial Z}{\partial b} = \frac{\partial Z}{\partial z} \cdot \frac{\partial z}{\partial b}$$

as $Z = w_1 x_1 + w_2 x_2 + b$

$$\boxed{\frac{\partial Z}{\partial b} = 1}$$

$$\therefore \boxed{db = a - y}$$

\therefore Gradient descent updates for logistic regression for a single training example.

$$\boxed{h_0(x) - y}$$

$$w_1 := w_1 - \alpha dw_1 = w_1 - \alpha (a - y) x_1 = w_1 - \alpha dz x_1$$

$$w_2 := w_2 - \alpha dw_2 = w_2 - \alpha (a - y) x_2 = w_2 - \alpha dz x_2$$

$$b := b - \alpha db = b - \alpha (a - y) = b - \alpha dz$$

$$\boxed{h_0(x) - y} \rightarrow \begin{matrix} \text{it is} \\ \text{simultaneous} \\ \text{add proof} \end{matrix}$$

* Take away message

* Use of computational graph makes the calculation of gradients/derivatives very efficient.

(5)

→ m examples

→

Gradient descent on m examples

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

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

$$d\omega_1^{(i)}, d\omega_2^{(i)}, db^{(i)}$$

$$d\omega_1 \Rightarrow \frac{\partial}{\partial \omega_1} J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \omega_1} \mathcal{L}(a^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)}) x_1^{(i)}$$

$$d\omega_2 = \frac{\partial}{\partial \omega_2} J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial \omega_2} (\mathcal{L}(a^{(i)}, y^{(i)})) = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)}) x_2^{(i)}$$

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

$$\frac{\partial L}{\partial}$$

$$z = \omega^T x + b$$

$$a = \sigma(z)$$

$$J = - \left[y \log a + (1-y) \log (1-a) \right]$$

$$dz = a - y$$

$$dw_1 = x_1 dz$$

$$dw_2 = x_2 dz$$

$$db = dz$$

// forward pass

// dimension / Backward pass

// backward pass

n = 2
if n > 2 loop for j in range(n) dw_j = x_j dz

$$J = m$$

$$dw_1 = m$$

$$dw_2 = m$$

$$db = m$$

$$dw_1 = m$$

$$dw_2 = m$$

$$db = m$$

$$dw_1 = m$$

$$dw_2 = m$$

$$db = m$$

$$dw_1 = m$$

$$dw_2 = m$$

$$db = m$$

$$dw_1 = m$$

$$dw_2 = m$$

$$db = m$$

$$dw_1 = m$$

$$dw_2 = m$$

$$db = m$$

$$dw_1 = m$$

$$dw_2 = m$$

$$db = m$$

$$\omega_1 := \omega_1 - \alpha dw_1$$

$$\omega_2 := \omega_2 - \alpha dw_2$$

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

// update weights

2 Weakness: 2 For loops → make code len efficient

In Deep Learning Era, we have bigger datasets, so need to do training w/o loop

→ Vectorization of training code help to remove explicit for loop.

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

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

w & x are large vector $w \in \mathbb{R}^{n_x}$

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

Non vectorized implementation:

$$Z = 0$$

for i in range(n_x)

$$Z += w[i] * x[i]$$

$$Z += b$$

same

Vectorized implementation

$$Z = \text{np.dot}(w^T, x) + b$$

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

$$(1, 2) (2, 1)$$

$$w_1 x_1 + w_2 x_2 + b$$

(single instruction multiple data)
 $\text{np.dot}()$ ^{numpy} ~~is~~ ^{uses} ~~fast~~ ^{SIMD} instructions both in CPU or GPU mode, which ultimately run the code very fast.
 uses SIMD instruction that do computation in parallel.

$$\begin{aligned} w &= w - \alpha dw \\ b &= b - \alpha db \end{aligned}$$

$$\begin{aligned} dw &= \frac{\partial L}{\partial w} \\ db &= \frac{\partial L}{\partial b} \end{aligned}$$

⑧
code / Dem: for vectorized result

$C = np.dot(a, b)$

VS

for i in range(1000000)
 $C += a[i] * b[i]$

import time

$a = np.random.rand(1000000)$

// uniform distribution

$b = np.random.rand(1000000)$

$t1c = time.time()$

$C = np.dot(a, b)$

$t2c = time.time()$

print(C)

print("vectorized version: " + str(1000 * (t2c - t1c)) + "ms") // 51.5ms

// let see Non vectorized version.

$c = 0$

$t1c = time.time()$

for i in range(1000000)
 $C += a[i] * b[i]$

$t2c = time.time()$

print(C)

print("for loop: " + str(1000 * (t2c - t1c)) + "ms") // 300ms

* Take away

use built function to run code faster as they are
single instructions

→ Whenever possible, avoid loops

$$U = A \underset{\text{matrix}}{V} \underset{\text{vector}}{v}$$

$$U = \sum_j A_{ij} V_j$$

$2 \times 3 \quad 3 \times 1$

$$\begin{bmatrix} \quad \end{bmatrix}_{3 \times 3} U \begin{bmatrix} \quad \end{bmatrix}_{3 \times 1} = \begin{bmatrix} \quad \end{bmatrix}_{3 \times 1}$$

Non vectorized Imp

$$U = \text{np.zeros}((n, 1))$$

```
for i in range(3):
    for j in range(3):
        U[i] += A[i][j] * V[j]
```

VS

vectorized Imp

$$U = \text{np.dot}(A, V) \quad // \text{it eliminates 2 for loops}$$

$$\Rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}_{3 \times 3} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} 1 \\ 4 \\ 7 \end{bmatrix}_{3 \times 1}$$

$A \quad V \quad U$

$$U = \text{np.dot}(A, V)$$

$$\begin{aligned} \omega &= \omega - \alpha d\omega \\ b &= b - \alpha db \end{aligned}$$

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

* Vectors & matrix valued functions

→ if we need to apply the exponential operation on every element of a matrix/vector

$$V = \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix} \Rightarrow u = \begin{bmatrix} e^{u_1} \\ \vdots \\ e^{u_n} \end{bmatrix}$$

→ Non vectorized / np:

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

for i in range(n):

$$u[i] = math.exp(V[i])$$

→ Vectorized / np:

• Numpy have many built in function that work in vectorized fashion.

$$\underline{u} = np.exp(V)$$

$$np.log(V)$$

$$np.abs(V)$$

$$np.maximum(V, 0)$$

$$V \times \times 2$$

// elementwise square

$$1/V$$

// elementwise inverse

→ Python operators do elementwise operation efficiently

1.1 Logistic regression algorithm to make loop free code.

First Attempt

$$J = 0, \quad dw_1 = 0, \quad \cancel{dw_2 = 0}, \quad db = 0$$

$$w = \begin{bmatrix} 1 \\ 1 \end{bmatrix}_{n_x}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n_x} \end{bmatrix}_{n_x}$$

for $i = 1$ to m :

accumulator $\hookrightarrow dw = \text{np.zeros}((n_x, 1))$

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

$$\hat{y}^{(i)} = \sigma(z^{(i)})$$

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

$$\textcircled{1} dz^{(i)} = y^{(i)} - \hat{y}^{(i)}$$

$$\text{for } j \text{ in range}(n_x): \quad \begin{cases} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_j += x_j^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \end{cases}$$

vector & scalar multiply

$$dw_1 += x_1 dz$$

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

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

$$dw /= m$$

Now we have one loop instead of 2 loops which make the program much faster.

$$\begin{cases} w = w - \alpha dw \\ b = b - \alpha db \end{cases}$$

Vectorized Logistic Regression (12)

* To speed up implementation of Training process — process the entire training set — using the idea of vectorization

Forward Pass

* currently forward propagation steps for m training examples inside a loop

$$\begin{array}{l|l|l} Z^{(1)} = W^T X^{(1)} + b & Z^{(2)} = W^T X^{(2)} + b & Z^{(3)} = W^T X^{(3)} + b \\ a^{(1)} = \sigma(Z^{(1)}) & a^{(2)} = \sigma(Z^{(2)}) & a^{(3)} = \sigma(Z^{(3)}) \end{array}$$

• w/o loop, how we perform computation?

$$X = \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & X^{(m)} \\ | & | & | \end{bmatrix} \in \mathbb{R}^{n_x, m}$$

(n_x, m)

$$W = \begin{bmatrix} | \\ | \\ | \end{bmatrix} \Rightarrow W^T = \begin{bmatrix} \text{---} \end{bmatrix} \quad \& \quad b = [b \ b \ \dots \ b]$$

(1, n_x) (1, m)

(n_x, 1)

$$W^T * \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & X^{(m)} \\ | & | & | \end{bmatrix} + [b \ b \ \dots \ b]$$

(1, n_x) (n_x, m) 1, m

broadcasting

$$= \begin{bmatrix} W^T X^{(1)} + b & W^T X^{(2)} + b & \dots & W^T X^{(m)} + b \end{bmatrix}$$

// stack of Z's for all training examples

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

// automatic

$$A = Z = \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} \Rightarrow A = \omega(Z)$$

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

$$A = \omega(Z)$$

$$A = \frac{1}{(1 + \text{np.exp}(-Z))}$$

vectorized soluti ✓

Backward Pass

• As we know

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

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

$$\therefore dz = \begin{bmatrix} dz^{(1)} & dz^{(2)} & \dots & dz^{(m)} \end{bmatrix} \quad // \text{stacking horizontally} \quad 1 \times m$$

$$A = \begin{bmatrix} a^{(1)} & \dots & a^{(m)} \end{bmatrix} \quad // \text{stack horizontally}$$

$$Y = \begin{bmatrix} y^{(1)} & \dots & y^{(m)} \end{bmatrix} \quad // \text{stack horizontally}$$

$$\checkmark \boxed{dz = A - Y} = \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & \dots & a^{(m)} - y^{(m)} \end{bmatrix} \quad (1 \times m)$$

$$\checkmark \boxed{dw = \frac{1}{m} X^T dz} \quad // \text{need to do transpose } dz \text{ as } dz \text{ is } (1, m) \text{ matrix}$$

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

G.D

$$\begin{cases} w = w - \alpha dw \\ b = b - \alpha db \end{cases}$$

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

Broad casting in python

- calories from carbs, proteins, Fats in 100 grams of different food.

	Apple	Beef	Eggs	potatos
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

$$= A$$

$\text{Total} = 59 \text{ cal in apple}$

- calculate % of calories from carb, protein, and fat.
- Can we do it without explicit for loop?

$$\text{cal} = A \cdot \text{sum}(\text{axis}=0) \quad // \text{vertically add.}$$

$$\text{Percentage} = 100 * A / \text{cal} \cdot \text{reshape}(1, 4)$$

(15)

More examples regarding broadcasting

$$(w^T x + b)$$

let us \Rightarrow

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

①

Python broadcasting operation
generates a $[4 \times 1]$ matrix as
shown below. Then performs the
Addition operation

$$= \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}_{4 \times 1} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix}_{4 \times 1} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}_{4 \times 1}$$

② Similarly

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

Python performs broadcasting operation & generates $[2 \times 3]$ matrix.
Corresponds to $[1 \times 3]$

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

$$= \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

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

(2,3) (2,1)

Python make 2x3 from 2,1 matrix by copying the column values 3 times.

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

$$= \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

• General principle for broadcasting

1) If (m,n) matrix & perform arithmetic operation with (1,n) matrix

$$\begin{bmatrix} + \\ / \\ * \\ / \end{bmatrix}$$

→ (m,n) broadcasting

2) If (m,n) matrix & $\begin{bmatrix} + \\ / \\ * \\ / \end{bmatrix}$ operation (m,1) matrix → (m,n) broadcast

3) If (m,1) & $\begin{bmatrix} + \\ / \\ * \\ / \end{bmatrix}$ with Real num → (m,1) by copy the real value
 e.g. $\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 \rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \end{bmatrix}$

4) (1,n) & $\begin{bmatrix} + \\ / \\ * \\ / \end{bmatrix}$ with Real num → (1,n)

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

• Tip to avoid bugs in code

```
import numpy as np
```

```
a = np.random.randn(5)
```

```
print(a)
```

↳ Gaussian variable

```
print(a.shape) // (5,) is Rank 1 array
```

```
print(a.T)
```

↳ not a row or a column vector

```
print(np.dot(a, a.T)) // is  $a \cdot a^T$ , but we may expect vector
```

Solu

```
a = np.random.randn(5, 1) // column vector
```

```
print(a.shape)
```

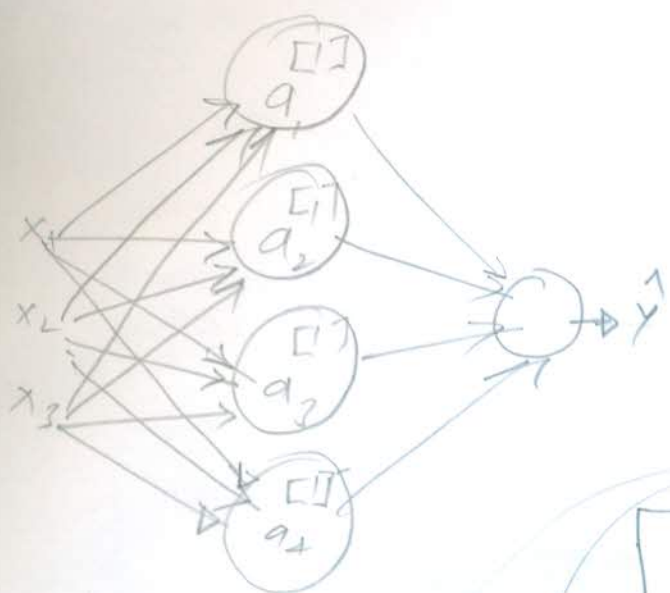
```
print(a.T) // row vector (1, 5)
```

5, 1, 1, 5

↳ (5, 1)

```
assert(a.shape == (5, 1)) // to check AssertionError
```


Neural NN representation



$$\begin{aligned}
 z_1 &= w_1 x + b_1, & a_1 &= \sigma(z_1) \\
 z_2 &= w_2 x + b_2, & a_2 &= \sigma(z_2) \\
 z_3 &= w_3 x + b_3, & a_3 &= \sigma(z_3) \\
 z_4 &= w_4 x + b_4, & a_4 &= \sigma(z_4)
 \end{aligned}$$

$$\begin{aligned}
 & \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}^T \text{ is a vector} \\
 & \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} \\
 & = \begin{bmatrix} w_1 x + b_1 \\ w_2 x + b_2 \\ w_3 x + b_3 \\ w_4 x + b_4 \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{bmatrix} = z
 \end{aligned}$$

⇒ Rule of Thumb: stack W & z vertically as nodes are placed vertically in net

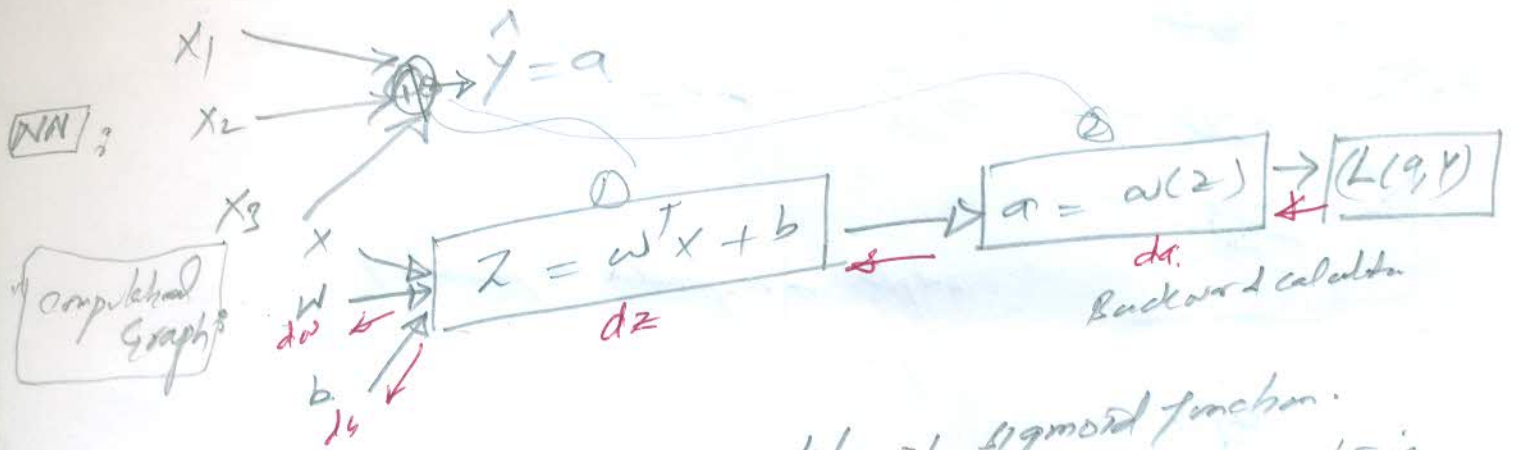
$$\sigma(z_1) = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix} = a$$

What is NN?

25/01/2011

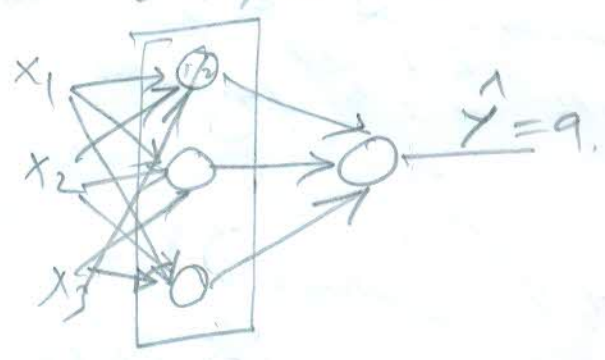
For logistic regression algorithm, the following computational graph we have used.

Forward Calculation

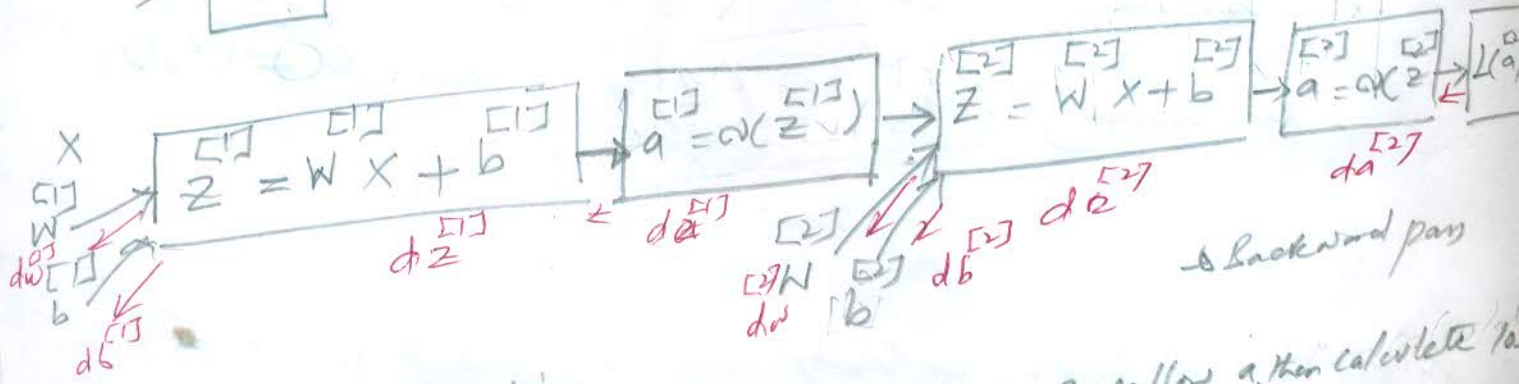


NN stacking together lot of sigmoid function.

- z followed by a calculation is done in multiple times
- logistic regression repeating twice.



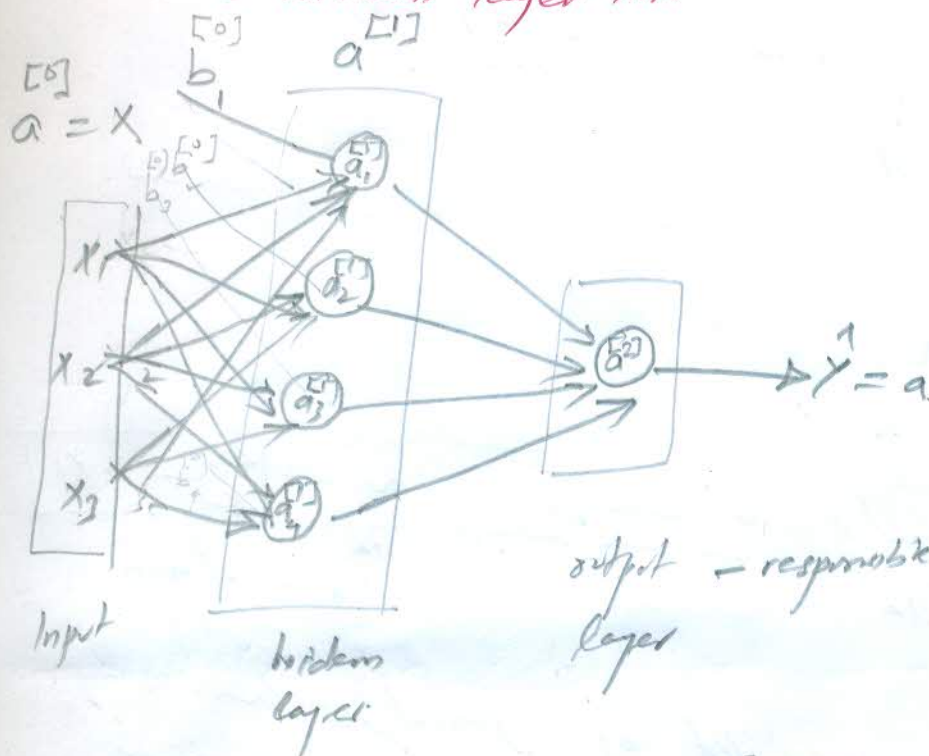
Forward pass



In NN, we have multiple z follow a , z follow a then calculate loss.

Similarly backward pass is to calculate derivatives.

- Neural Network presentation : Details
- One hidden layer NN



output - responsible to produce \hat{y}
layer

- what are calculation $a^{[0]} = x$, $a^{[1]}$, $a^{[2]} = y$

- it is 2 layer NN, not counting input layer
- hidden layer associated parameters $w^{[1]}$, $b^{[1]}$, whereas output layer is associated with $w^{[2]}$, $b^{[2]}$ parameters

- Dimension of $w^{[1]}$ is $(4, 3)$, where $b^{[1]} = (4, 1)$

- Dimension of $w^{[2]}$ is $(1, 4)$, where $b^{[2]} = (1, 1)$

1 unit 4 unit

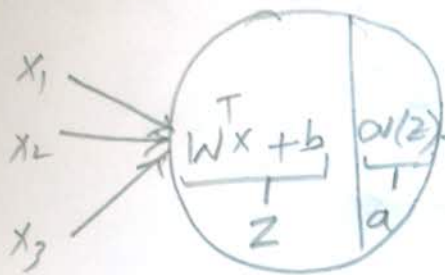
- What is the network actually computing, see it Next Page.

27- Detail of what NN with one hidden layer is computing:

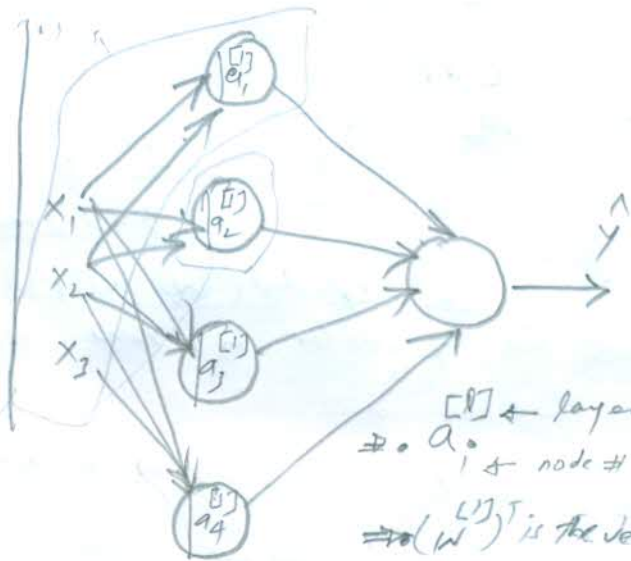
→ In logistic Regression a single unit compute two things ① Z followed by ② a

$$Z = W^T X + b$$

$$a = \sigma(Z)$$



$$a = \hat{y}$$



$[l]$ ← layer
 $\neq a_i$ ← node # in a layer
 $\Rightarrow (W^{[l]})^T$ is the vector transpose

$$Z_1^{[l]} = W_1^{[l]T} X + b_1^{[l]}$$

$$a_1^{[l]} = \sigma(Z_1^{[l]})$$

$$Z_2^{[l]} = W_2^{[l]T} X + b_2^{[l]}$$

$$a_2^{[l]} = \sigma(Z_2^{[l]})$$

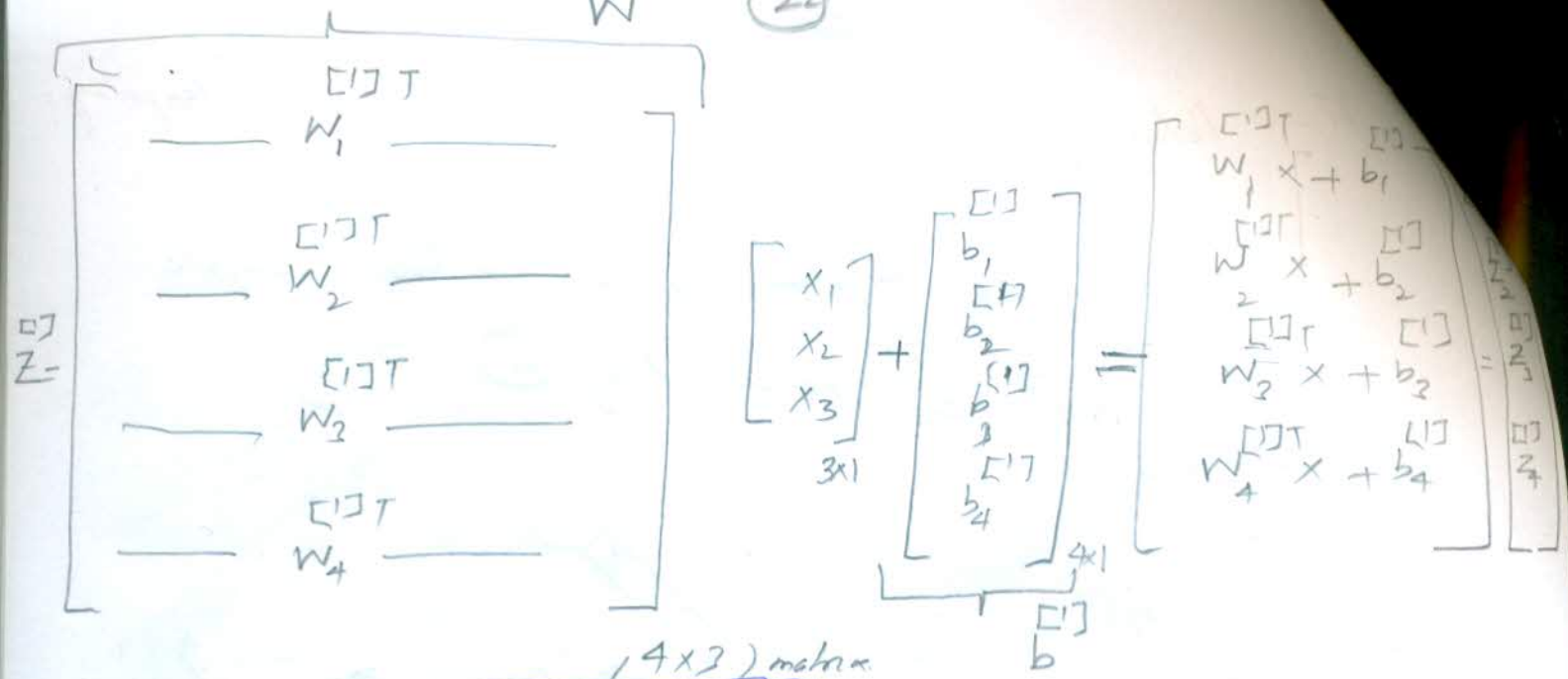
$$Z_3^{[l]} = W_3^{[l]T} X + b_3^{[l]}$$

$$a_3^{[l]} = \sigma(Z_3^{[l]})$$

$$Z_4^{[l]} = W_4^{[l]T} X + b_4^{[l]}$$

$$a_4^{[l]} = \sigma(Z_4^{[l]})$$

Do calculation using for loop is inefficient, so take & evaluate 2 Vectorized them.



$$z = \text{np.dot}(W, x) + b$$

Stack the 4 vector in a row as shown above.

~~1~~ Vectorization Rule of Thumb, when we have different node in a layer, we stack them vertically - so similarly we stack their results vertically in matrix

we call $(4, 3)$ weight matrix as W and $(4, 1)$ bias matrix as b at hidden unit

Similarly

$$a = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix} = \text{act}(z)$$

↑
element

$$A = \text{np.tanh}(z)$$

(22)

Given input x

$$z_1 = \text{np.dot}(w_1, x) + b_1$$

$$a_1 = \text{np.tanh}(z_1)$$

$$z_2 = \text{np.dot}(w, a_1) + b_2$$

$$a_2 = \text{sigmoid}(z_2)$$

$$z^{[1]} = w^{[1]} x^{[1]} + b^{[1]}$$

$$(4,1) = (4,3) (3,1) + (4,1)$$

$$a^{[1]} = \sigma(z^{[1]})$$

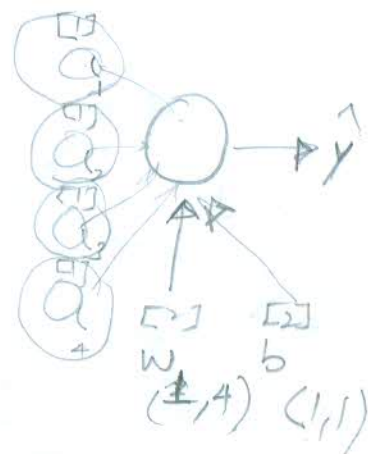
$$(4,1)$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$(1,1) = (1,4) (4,1) + (1,1)$$

$$a^{[2]} = \sigma(z^{[2]})$$

$$(1,1)$$



$$W^T = W$$

$$b = b$$

* If we just think the output unit as the logistic regression unit

$$W^{[2]} = W^T$$

$$b^{[2]} = b$$

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

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

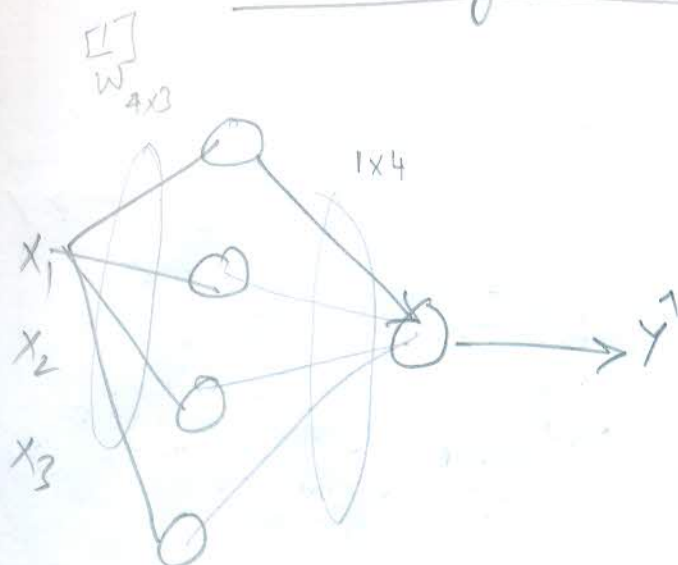
* We have to compute only above 4 equations

for computing \hat{y} for a single training example.

* now let see vectorized: across multiple examples

ONE hidden layer Neural Network

Vectorizing across multiple examples



$$Z^{[1]} = W^{[1]} x + b^{[1]}$$

$$a^{[1]} = \sigma(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(Z^{[2]})$$

$$X \longrightarrow a^{[2]} = \hat{y}$$

$$x^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)}$$

$$x^{(2)} \longrightarrow a^{2} = \hat{y}^{(2)}$$

$$\vdots$$

$$x^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)}$$

(for a single Training examples)

$a^{[2](i)}$ example i
layer 2

~~XXX~~

For $i = 1$ to m :

$$Z^{[1](i)} = W^{[1]} x^{(i)} + b^{[1]}$$

$(4 \times 3) \quad (3 \times 1)$

$$a^{[1](i)} = \sigma(Z^{[1](i)})$$

(4×1)

$$Z^{[2](i)} = W^{[2]} a^{[1](i)} + b^{[2]}$$

$(1 \times 4) \quad (4 \times 1)$

$$a^{[2](i)} = \sigma(Z^{[2](i)})$$

* We can also avoid this for loop also using vectorization.

$$\begin{bmatrix} Z^{1} \\ \vdots \\ Z^{[1](m)} \end{bmatrix} = Z$$

vectorized across multiple

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

3 100

per batch, we have to compute the following

$$Z^{(1)} = W^{(0)} X + b^{(1)} \quad (\text{if } m=100)$$

$$A^{(1)} = \sigma(Z^{(1)})$$

$$Z^{(2)} = W^{(1)} A^{(1)} + b^{(2)}$$

$$A^{(2)} = \sigma(Z^{(2)})$$

$$\downarrow \rightarrow Z^{(3)} =$$

$$\begin{bmatrix} Z_1^{(1)}(1) \\ Z_2^{(1)}(1) \\ Z_3^{(1)}(1) \\ Z_4^{(1)}(1) \end{bmatrix} \quad \begin{bmatrix} Z_1^{(1)}(2) \\ Z_2^{(1)}(2) \\ Z_3^{(1)}(2) \\ Z_4^{(1)}(2) \end{bmatrix} \quad \dots$$

previously, we have sample

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

$$Z^{(1)}(i) = W^{(0)} x^{(i)} + b^{(1)}$$

$$a^{(1)}(i) = \sigma(Z^{(1)}(i))$$

$$Z^{(2)}(i) = W^{(1)} a^{(1)}(i) + b^{(2)}$$

$$a^{(2)}(i) = \sigma(Z^{(2)}(i))$$

$$\begin{bmatrix} Z_1^{(1)}(m) \\ Z_2^{(1)}(m) \\ Z_3^{(1)}(m) \\ Z_4^{(1)}(m) \end{bmatrix}$$

($Z^{(1)}$ obtained by stacking $Z^{(1)}(i)$ horizontally)

$$\begin{bmatrix} a_1^{(1)}(m) \\ a_2^{(1)}(m) \\ a_3^{(1)}(m) \\ a_4^{(1)}(m) \end{bmatrix}$$

($a^{(1)}$ obtained by stacking $a^{(1)}(i)$ horizontally)

training examples A
hidden units at hidden unit

$$A = \begin{bmatrix} a_1^{(1)}(1) & a_1^{(1)}(2) & \dots \\ a_2^{(1)}(1) & a_2^{(1)}(2) & \dots \\ a_3^{(1)}(1) & a_3^{(1)}(2) & \dots \\ a_4^{(1)}(1) & a_4^{(1)}(2) & \dots \end{bmatrix}$$

$$A1 = \sigma_p \cdot \tanh(Z1)$$

⑤

②⑥

Justification for correction of vectorized implementation

Let consider 3 training examples

$$Z = W X + b$$

For 1st training example

$$Z^{(1)} = W X^{(1)} + b$$

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

For 2nd training example

$$Z^{(2)} = W X^{(2)} + b$$

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

For 3rd training example

$$Z^{(3)} = W X^{(3)} + b$$

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

①

$$W = \begin{bmatrix} \quad \\ \quad \\ \quad \\ \quad \end{bmatrix}_{4 \times 3}$$

$$W X^{(1)} = \begin{bmatrix} \quad \\ \quad \\ \quad \\ \quad \end{bmatrix}_{4 \times 1} \begin{bmatrix} \quad \\ \quad \\ \quad \\ \quad \end{bmatrix}_{3 \times 1} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}_{4 \times 1}$$

②

$$X = \begin{bmatrix} X^{(1)} & X^{(2)} & X^{(3)} \\ x_1^{(1)} & x_1^{(2)} & x_1^{(3)} \\ x_2^{(1)} & x_2^{(2)} & x_2^{(3)} \\ x_3^{(1)} & x_3^{(2)} & x_3^{(3)} \end{bmatrix}_{3 \times 3}$$

$$W X^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}_{4 \times 1}$$

$$W X^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}_{4 \times 1}$$

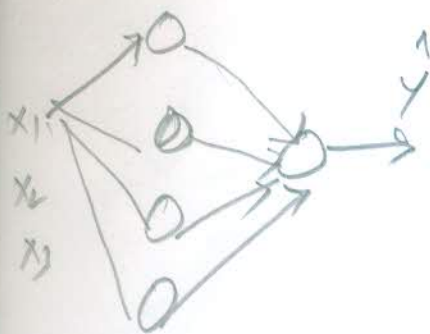
$$W \begin{bmatrix} X^{(1)} & X^{(2)} & X^{(3)} \\ x_1^{(1)} & x_1^{(2)} & x_1^{(3)} \\ x_2^{(1)} & x_2^{(2)} & x_2^{(3)} \\ x_3^{(1)} & x_3^{(2)} & x_3^{(3)} \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}_{4 \times 3} = \begin{bmatrix} Z^{(1)} & Z^{(2)} & Z^{(3)} \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}_{4 \times 3} = Z$$

\vdots bias
 b bias
 b bias
 b bias

we can add b using broadcasting

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

Recap of vectorizing across multiple examples.



for $i = 1$ to m :

$$Z^{[1]}(i) = W^{[1]} x^{(i)} + b^{[1]}$$

$$a^{[1]}(i) = \sigma(Z^{[1]}(i))$$

$$Z^{[2]}(i) = W^{[2]} a^{[1]}(i) + b^{[2]}$$

$$a^{[2]}(i) = \sigma(Z^{[2]}(i))$$

Show correctness of following the following.

let

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(m)} \\ | & | & | \end{bmatrix}_{3 \times m}$$

$$A = \begin{bmatrix} | & | & | \\ a^{[1]}(1) & a^{[1]}(2) & a^{[1]}(m) \\ | & | & | \end{bmatrix}_{4 \times m}$$

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]})$$

⇒ Gradient Descent for Neural Network with single hidden unit

Parameters & dimensions

$$\begin{array}{cc} \begin{array}{c} [1] \\ W, b \\ (n^{[0]}, n^{[1]}) \end{array} & \begin{array}{c} [1] \\ W, b \\ (n^{[1]}, n^{[2]}) \end{array} \\ & \begin{array}{c} [2] \\ W, b \\ (n^{[2]}, n^{[3]}) \end{array} \end{array}$$

- $n^{[0]}$ are no input features

- $n^{[1]}$ are no of hidden units at layer [1]

- $n^{[2]}$ are the no of hidden units at layer [2]

Cost function: for binary classification

$$J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^i, y^i)$$

Gradient Descent: // Initialize the parameters randomly, instead of initialize them with zeros

Repeat {
forward pass. ~~to compute~~ ① compute prediction ($\hat{y}^{(i)}$; $i=1, \dots, m$)
compute derivatives ② $\frac{dJ}{dW^{[1]}} = \frac{dJ}{dW^{[1]}}$, $\frac{dJ}{db^{[1]}} = \frac{dJ}{db^{[1]}} \dots$

③ update parameters

$$W^{[1]} = W^{[1]} - \alpha \frac{dJ}{dW^{[1]}}$$

$$b^{[1]} = b^{[1]} - \alpha \frac{dJ}{db^{[1]}}$$

$$W^{[2]} = W^{[2]} - \alpha \frac{dJ}{dW^{[2]}}$$

$$b^{[2]} = b^{[2]} - \alpha \frac{dJ}{db^{[2]}}$$

}

⇒ Equations for Forward & backward propagation

• Forward pass // all are vectorised

$$Z^{[1]} = W^{[0]} X + b^{[0]}$$

$$A^{[1]} = g(Z^{[1]})$$

$$Z^{[2]} = W^{[1]} A^{[1]} + b^{[1]}$$

as it doesn't change problem

$$A^{[2]} = g(Z^{[2]}) = \sigma(Z^{[2]})$$

• Backward pass to calculate derivatives

let $Y = \{y^{(1)}, y^{(2)} \dots y^{(m)}\}$ // Grand Truth

$$dZ^{[2]} = \frac{dL}{dZ^{[2]}} = \frac{1}{m} A^{[2]'} - Y$$

$(n^{[2]}, 1)$

$(n^{[2]}, 1)$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis=1, keepdims=True)$$

Similar to logistic Regression

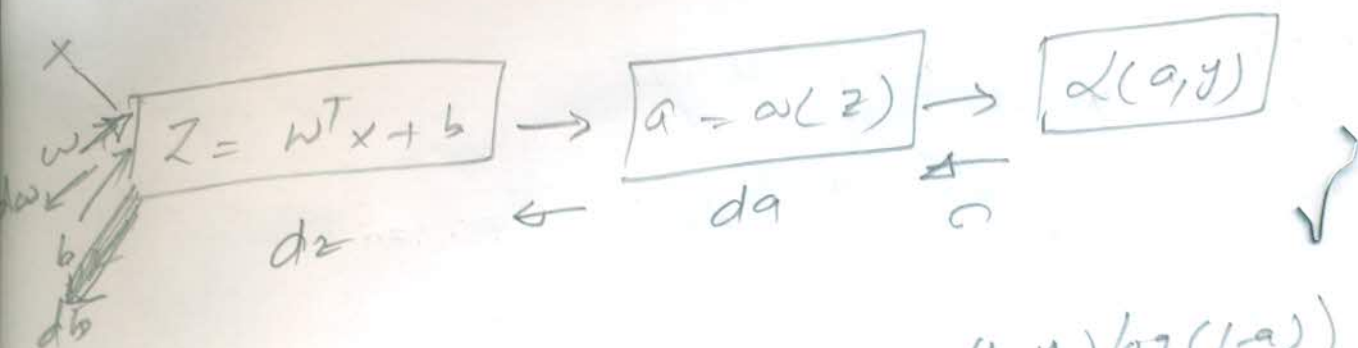
$$dZ^{[1]} = \frac{1}{m} dZ^{[2]} W^{[1]T} * g'(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis=1, keepdims=True)$$

$(n^{[1]}, 1)$
 $(n^{[1]}, 1)$

Computational Graph for Logistic Regression



As $L(a, y) = -(y \log a + (1-y) \log(1-a))$

$da = \frac{\partial L(a, y)}{\partial a} = -\left(\frac{y}{a} + \frac{1-y}{1-a}\right) = \left[\frac{-y}{a} + \frac{1-y}{1-a}\right]$ ✓

As $dz = \frac{\partial L}{\partial z} = \left[\frac{\partial L}{\partial a} * \frac{\partial a}{\partial z}\right] = da \cdot \frac{\partial a}{\partial z}$ (Chain rule of calculus)

As for logistic regression we use sigmoid for $a = \frac{1}{1 + e^{-z}}$

$\therefore \frac{\partial a}{\partial z} = a(1-a) = \frac{1}{(1+e^{-z})} \cdot \left(1 - \frac{1}{1+e^{-z}}\right)$

$\therefore dz = \left(\frac{-y}{a} + \frac{1-y}{1-a}\right) a(1-a) = \cancel{a} \cdot \cancel{y} + a - \cancel{a} \cdot \cancel{y} = a - y$

✓ $\boxed{dz = a - y}$ ✓

$dw = \frac{\partial L(a, y)}{\partial w} = \frac{\partial L(a, y)}{\partial z} \cdot \frac{\partial z}{\partial w} = dz \cdot \frac{\partial z}{\partial w}$

As $z = w^T x + b \Rightarrow \boxed{\frac{\partial z}{\partial w} = x}$

✓ $\boxed{dw = dz \cdot x}$

$db = \frac{\partial L(a, y)}{\partial b} = \frac{\partial L(a, y)}{\partial z} \cdot \frac{\partial z}{\partial b} = dz \cdot \frac{\partial z}{\partial b}$

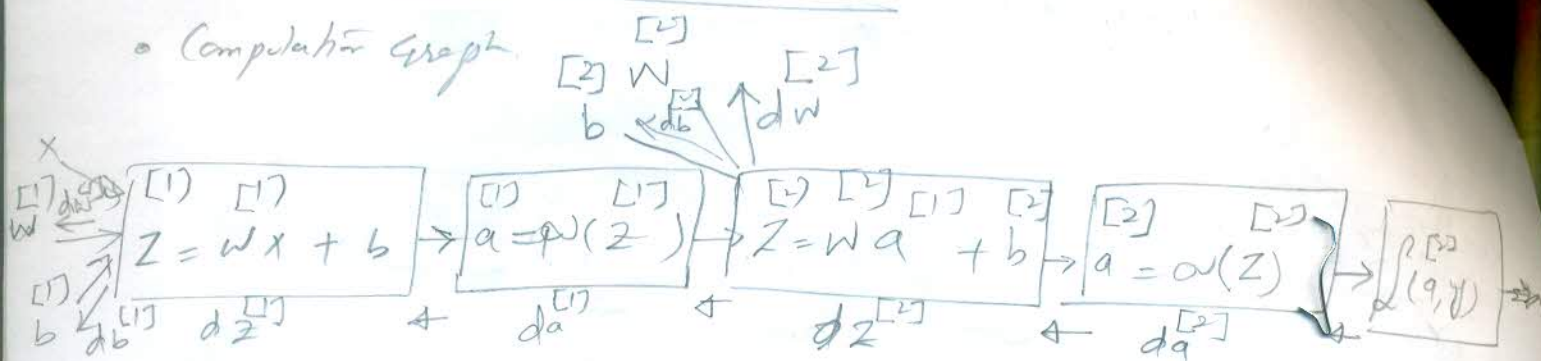
$\boxed{\frac{\partial z}{\partial b} = 1}$

✓ $\boxed{db = dz}$

(31)

Gradients with 2 layer NN

• Computation Graph



$$da^{[2]} = \frac{\partial L(a, y)}{\partial a^{[2]}} = \frac{\partial}{\partial a^{[2]}} \left[-(y \log a^{[2]} + (1-y) \log (1-a^{[2]})) \right]$$

$$= -\frac{y}{a^{[2]}} - (1-y) \cdot \frac{1}{1-a^{[2]}} \cdot (-1)$$

$$da^{[2]} = -\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}}$$

$$dz^{[2]} = \frac{\partial L(a, y)}{\partial z^{[2]}} = \frac{\partial L(a, y)}{\partial a^{[2]}} * \frac{\partial a^{[2]}}{\partial z^{[2]}} = da^{[2]} * \frac{\partial a^{[2]}}{\partial z^{[2]}}$$

$$a^{[2]} = \frac{1}{1 + e^{-z^{[2]}}}$$

$$\frac{\partial a^{[2]}}{\partial z^{[2]}} = a^{[2]} (1 - a^{[2]})$$

$$\rightarrow dz^{[2]} = \left(-\frac{y}{a^{[2]}} + \frac{1-y}{1-a^{[2]}} \right) * a^{[2]} (1 - a^{[2]})$$

$$dz^{[2]} = -y + y a^{[2]} + a^{[2]} - y a^{[2]} = a^{[2]} - y$$

$$dz^{[2]} = a^{[2]} - y$$

$$\frac{\partial L}{\partial w^{[2]}} = \frac{\partial L}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial w^{[2]}} = dz^{[2]} * \frac{\partial z^{[2]}}{\partial w^{[2]}}$$

Dimensions:
 $\frac{\partial L}{\partial z^{[2]}}$ [1]
 $\frac{\partial z^{[2]}}{\partial w^{[2]}}$ [2]
 $z^{[2]}, w^{[2]}, b^{[2]}$

as

$$z = w a + b$$

$$\frac{\partial z^{[2]}}{\partial w^{[2]}} = a^{[1]}$$

$$dw^{[2]} = \frac{\partial L}{\partial w^{[2]}} = (a - y)^{[2]} \times a^{[1]T}$$

$$dw^{[2]} = (a - y)^{[2]} \times a^{[1]T}$$

$$\frac{\partial L}{\partial b^{[2]}} = \frac{\partial L}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial b^{[2]}} = dz^{[2]} \times \frac{\partial z^{[2]}}{\partial b^{[2]}}$$

$$\frac{\partial z^{[2]}}{\partial b^{[2]}} = 1$$

$$db^{[2]} = \frac{\partial L}{\partial b^{[2]}} = dz^{[2]} \times 1 = dz^{[2]}$$

$$db^{[2]} = dz^{[2]}$$

$$da^{[1]} = \frac{\partial L}{\partial a^{[1]}} = \frac{\partial L}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} = dz^{[2]} \times \frac{\partial z^{[2]}}{\partial a^{[1]}}$$

As $z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$

$$\frac{\partial z^{[2]}}{\partial a^{[1]}} = w^{[2]}$$

$$da^{[1]} = dz^{[2]} \times w^{[2]} = (a - y)^{[2]} w^{[2]}$$

Previously:

$$dw = dz \times x$$

$$w^{[1]} = []$$

dim?

w is row vector so need Transpose in NN.
But in logistic Reym no Transpose is required
→ dim of w & dw are same.

$$N^{[2]} = [\dots]$$

$$dw^{[2]} = [\dots]$$

need Transpose for vector
not need for a single value.

where in case of $\frac{\partial L}{\partial w}$
 dw is Coln vector
so no need of Transpose.

$$as \quad z = W a + b$$

$$\frac{\partial z}{\partial W^{[2]}_{[1]}} = a^{[1]}$$

$$\therefore dW = \frac{\partial L}{\partial W^{[2]}_{[1]}} = (a - y) \times a^{[1]T}$$

$$dW = (a - y) \times a^{[1]T}$$

Previously:

$$\Rightarrow dW = dz \cdot x$$

$$W = [\quad]^{[1]}_{[2]}$$

dim?

W is row vector so need Transpose in NN.
But in log prob Regm no transpose is required
→ dim of W & dW are same

$$\therefore \frac{\partial L}{\partial b^{[2]}} = \left(\frac{\partial L}{\partial z^{[2]}} \right) \times \frac{\partial z^{[2]}}{\partial b^{[2]}} = dz \times \frac{\partial z^{[2]}}{\partial b^{[2]}}$$

$$\Rightarrow \frac{\partial z^{[2]}}{\partial b^{[2]}} = 1$$

$$db = \frac{\partial L}{\partial b^{[2]}} = dz \times 1 = dz$$

$$db = dz$$

W = [- - -]
dW = [- - -]
need Transpose for vector
not need for a single value.

where in case of
 $dW = dz \cdot x$
dW is Coln vector
so no need of Transpose.

$$da = \frac{\partial L}{\partial a^{[1]}} = \frac{\partial L}{\partial z^{[2]}} \times \frac{\partial z^{[2]}}{\partial a^{[1]}} = dz \times \frac{\partial z^{[2]}}{\partial a^{[1]}}$$

$$As \quad z = W a + b$$

$$\therefore \frac{\partial z^{[2]}}{\partial a^{[1]}} = W$$

$$\therefore da = dz \times W = (a - y) W$$

$$da = (a - y) W$$

$$\frac{\partial L}{\partial z^{[l]}} = \frac{\partial L}{\partial a^{[l]}} \cdot \frac{\partial a^{[l]}}{\partial z^{[l]}} = da^{[l]} \cdot \frac{\partial a^{[l]}}{\partial z^{[l]}}$$

$$a^{[l]} = g(z^{[l]})$$

Any Activation Function

$$\frac{\partial a^{[l]}}{\partial z^{[l]}} = g'(z^{[l]})$$

$$dz^{[l]} = da^{[l]} \cdot g'(z^{[l]})$$

$$dz^{[l]} = dz^{[l-1]} w^{[l-1]} \cdot g'(z^{[l]})$$

$$dz^{[l]} = w^{[l-1]T} dz^{[l-1]} \cdot g'(z^{[l]})$$

$$\frac{\partial L}{\partial w^{[l]}} = \frac{\partial L}{\partial z^{[l]}} \cdot \frac{\partial z^{[l]}}{\partial w^{[l]}} = dz^{[l]} \cdot \frac{\partial z^{[l]}}{\partial w^{[l]}}$$

$$\text{As } z^{[l]} = w^{[l]} x + b^{[l]}$$

$$\frac{\partial z^{[l]}}{\partial w^{[l]}} = x$$

$$\therefore dw^{[l]} = dz^{[l]} \cdot x^T$$

dim?

$$db^{[l]} = dz^{[l]}$$