

5 *unseen*

10 *beginning*

15 *feel*

20

25

30

what I will Learn

1. Neural Networks & Deep learning

5 2. Improving Deep Neural Networks:-

3. Structuring ML Project

10 4. CNN | 8. NLP: Sequence Models.

RNN / LSTM

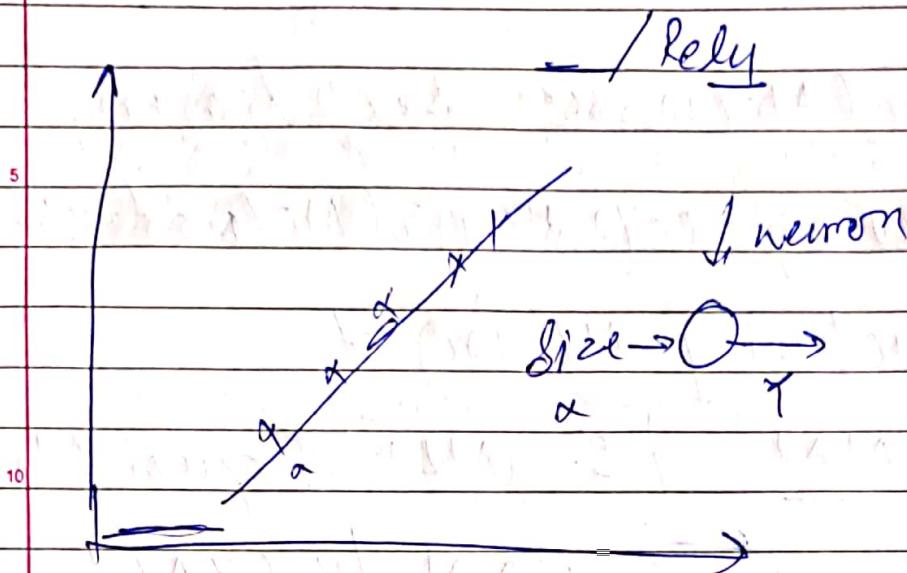
15

20

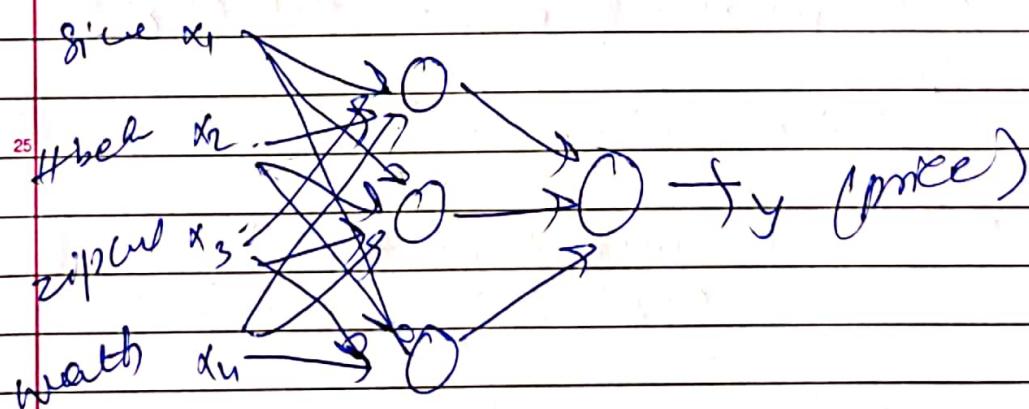
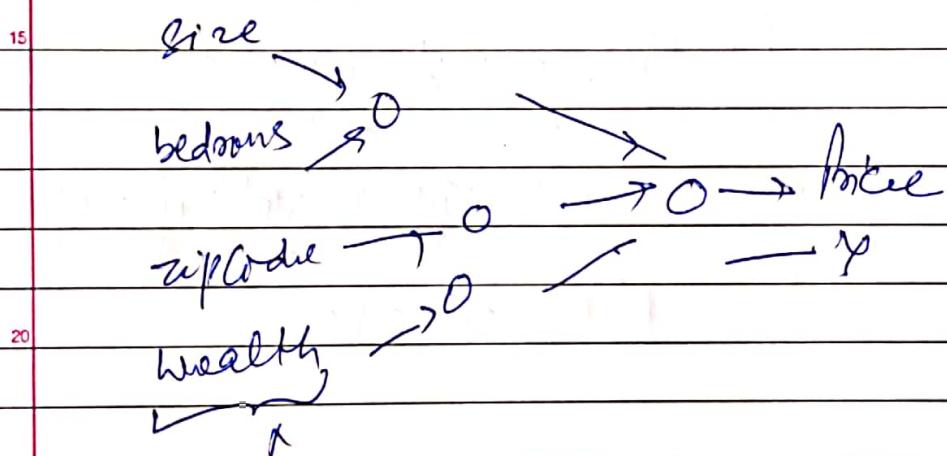
25

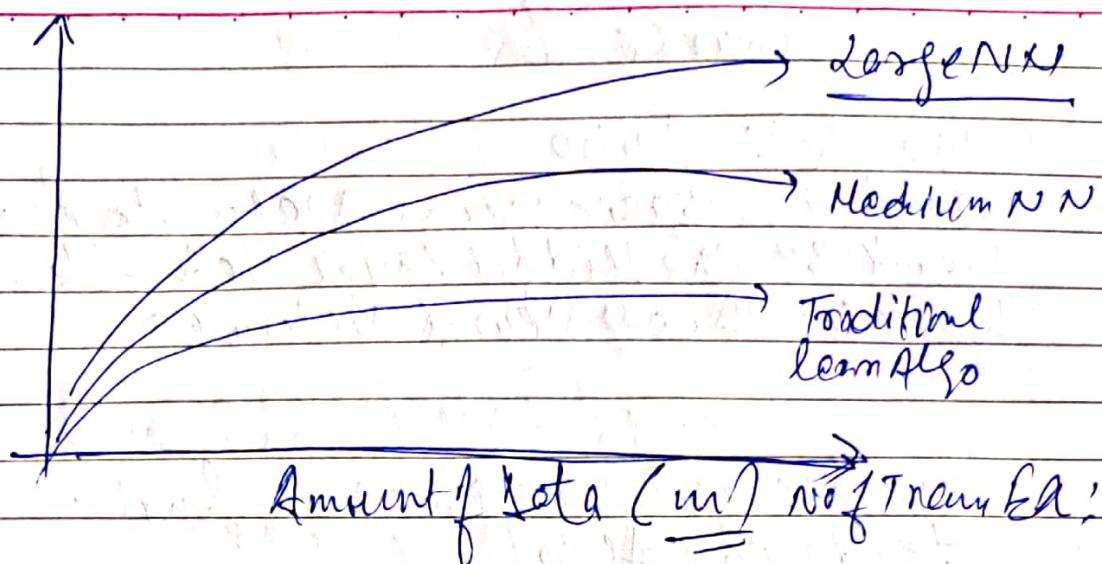
30

Neural Network



Housing Price Prediction





Small Training Set (SVM) can do better
whereas Neural Network requires
by more f more data.

Idea

Expert Code

Stochastic gradient descent in optimization

$(w_i \rightarrow w_i + \eta \cdot \delta_i)$

Local optima

Course: 6

Week 1: Intro

Week 2: Basic neural Network theory.

Week 3: One hidden layer Network.

Week 4: Deep Neural Networks.

Increase train size doesn't hurt
traditional Algo.

Increase size of neural network does
not hurt but it can help.

Week 2

Logistic Regression as $\sigma(\theta^T x)$

Neural Network

Notation - Single Train example

is denoted $(x, y) \in \mathbb{R}^n$

where x is n dimensional feature
vector. $y \in \{0, 1\}$

m (no of train examples)

$$(x^{(1)}, y^{(1)}) \quad (x^{(2)}, y^{(2)}) \quad \dots \quad (x^{(m)}, y^{(m)})$$

$M = M_{\text{train}}$ / No. of train example
 M_{test} / No. of test example

5 $X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}$ taking each train example
 stacking up into col.
 10 notice that in OpenCV
 we do it ~~as~~ \Rightarrow ROW wise stacking.

15 $\rightarrow X \cdot \text{shape}(n_x, m)$

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

Stacking labels into cols :-

20 $y \cdot \text{shape}(1, m)$

L. logistic Regression

25 Given X , what is \hat{y} (i.e.) label

$$\hat{y} = P(Y=1) \quad P(Y=1/X)$$

Param: $w \in \mathbb{R}^n$, $b \in \mathbb{R}$

30 Output: $\hat{y} = \sigma(w^T x + b) \rightarrow z$
 $\sigma = \text{sigmoid function}$

Logistic Regression Cost Function

$$y = \sigma(z) \text{ where } z = w^T x + b$$

Given $(x^{(0)}, y^{(0)}), (x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$
 want $\hat{y}^{(i)} \approx y^{(i)}$

$$d(\hat{y}, y) = \text{cross entropy loss}$$

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

if $y=1$

$$\Rightarrow d(\hat{y}, y) = -\log \hat{y}$$

if $y=0$

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

Cost Function

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m d(\hat{y}^{(i)}, y^{(i)})$$

$$= \frac{1}{m} \sum_{i=1}^m (-\hat{y}^{(i)} \log \hat{y}^{(i)} - (1-\hat{y}^{(i)}) \log(1-\hat{y}^{(i)}))$$

30

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

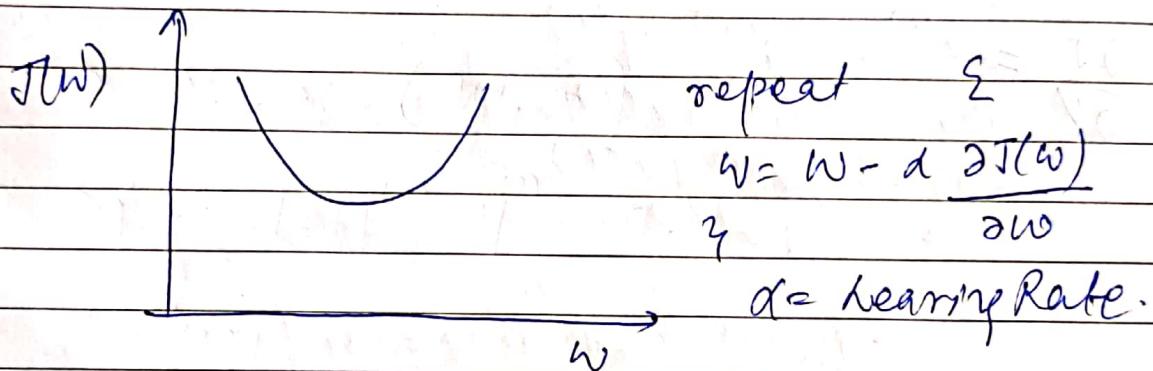
Gradient Descent

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m d(y^{(i)}, \hat{y}^{(i)})$$

$$= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})$$

want to find w, b which min $J(w, b)$

Usually for "logistic Regression" we init by "25%".



will use convention $dw = \frac{\partial J(w)}{\partial w}$

$$so \quad w = w - 1 dw \quad | \quad so$$

$$w = w - \alpha \left[\frac{\partial J}{\partial w} \right] \rightarrow dw$$

$$b = b - \alpha \left[\frac{\partial J}{\partial b} \right] \rightarrow db$$

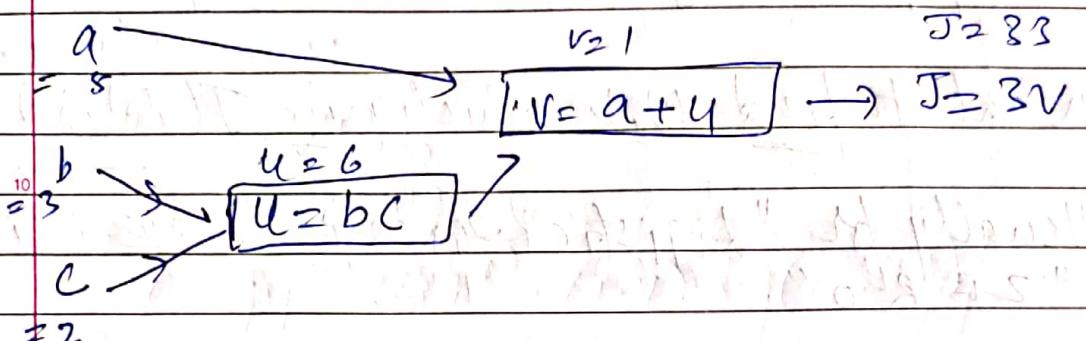
Computation Graph

$$J(a, b, c) = 3(a + bc)$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



$$\frac{\partial J}{\partial v} =$$

20

25

30

Logistic Regression Gradient Descent

$$\textcircled{1} \quad z = w_1 x_1 + b \quad | \quad y = a = \sigma(z)$$

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

Let say we have feature length 2' (x_1, x_2)

$$\begin{array}{l} x_1 \\ w_1 \\ x_2 \\ w_2 \\ b \end{array} \rightarrow \boxed{z = w_1 x_1 + w_2 x_2 + b} \rightarrow \boxed{\hat{y} = a = \sigma(z)} \quad | \quad L(a, y)$$

$$dz = \frac{\partial L(a, y)}{\partial z} \quad "da" = \frac{\partial a}{\partial z}$$

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

$$\frac{\partial L}{\partial a} \times \boxed{\frac{\partial a}{\partial z}} = -a(1-a)$$

$$\frac{\partial L(a, y)}{\partial w_1} = "dw_1" = x_1 \cdot dz \quad | \quad dw_2 = x_2 \cdot dz \quad | \quad db = dz$$

$$w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

Logistic regression On "m" examples

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

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

for $i=1 \rightarrow m$

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

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

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

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

$$dw_1 += x_1^{(i)} dz^{(i)} \quad | \quad n=2$$

$$dw_2 += x_2^{(i)} dz^{(i)} \quad | \quad (B-A) =$$

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

$$J \leftarrow J + \text{ } | \quad dw_1 = m |$$

$$w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

25

30

Logistic Regression Vectorization

Vectorization of code enable SIMD instruction.

$$5 \quad X = \begin{bmatrix} 1 & x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & 1 & \dots & 1 \end{bmatrix} \quad (n_x, m)$$

$$10 \quad [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = w^T X + [b, b_2 - b_1] \quad 1 \times m$$

$$w^T \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix} \quad (n_x, m) \quad =$$

15 w^T is a row vector

$$\begin{bmatrix} 1 & 1 & 1 \\ x_1 & x_2 & \dots & x_m \\ b & b & \dots & b \end{bmatrix}$$

$$20 \quad = [w^T x_1 \ w^T x_2 \ \dots \ w^T x_m] + [b, b_2 - b_1 \ \dots \ b_m]$$

$$\mathcal{Z} = [w^T x_1 + b, w^T x_2 + b_2, \dots, w^T x_m + b_m]$$

$$= [z_1 \ z_2 \ \dots \ z_m]$$

$$25 \quad \mathcal{Z} = np. \text{dot}(w^T, X) + b \quad \text{"Broadcasting"}$$

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

$$dZ^u = a^{(1)} - y^{(1)} \quad | \quad dZ^{(2)} = a^{(2)} - y^{(2)} \\ dZ = [dZ^{(1)} \quad dZ^{(2)} \quad \dots \quad dZ^{(m)}] \quad \underline{\text{1xu}}$$

$$A = [a^{(1)} \quad a^{(2)} \quad \dots \quad a^{(n)}] \cdot Y = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}]$$

5

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

$$dw = 0$$

10

$$dw^+ = X^{(1)} dZ^{(1)} \quad - \text{first train example}$$

$$dw^+ = X^{(2)} dZ^{(2)} \quad - \text{second train example}$$

15

$$dw^+ = X^{(m)} dZ^{(m)} \quad - \text{mth train example}$$

$$dw = dw/m$$

20

$$db^+ = dZ^{(1)} \quad - \text{1st train exam}$$

$$db^+ = dZ^{(2)} \quad - \text{2nd train example}$$

$$db = db/m$$

25

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

$$= \frac{1}{m} \cancel{np \cdot \text{sum}}(dZ)$$

30

$$dw = \frac{1}{m} X dZ^T$$

$$= \frac{1}{m} \left[\sum_{i=1}^m x^{(i)}_1 x^{(i)}_2 - \bar{x}^{(m)}_1 \bar{x}^{(m)}_2 \right] \begin{bmatrix} d\bar{x}^{(1)} \\ d\bar{x}^{(2)} \\ \vdots \\ d\bar{x}^{(m)} \end{bmatrix}$$

$$(n_{x,m}) (n_{m,1}) = (\underline{n_x \ 1})$$

$$= \frac{1}{m} \left[x_1 d\gamma + x^{(m)}_1 d\bar{x}^{(2)} - \bar{x}^{(m)}_1 d\bar{x}^{(m)} \right]$$

NonVectorize

$$\delta_2 = 0, dw_1 = 0, dw_2 = 0, db = 0$$

for $i=1 \rightarrow m$

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

Vectorize

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

$$= np \cdot \text{dot}(W.T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$J_T = - \left[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)}) \right] \quad dw = \frac{1}{m} X dZ^T$$

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

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

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

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

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

$$w = w - \lambda dw$$

$$b = b - \lambda db$$

$$J = J/m \quad dw_1 = dw_1/m$$

single
iteration gradient

descent

$$dw_2 = dw_2/m$$

$$db = db/m$$

30

of Broad Casting Python

	apple	Beef	egg	Potato
Carb	86.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	135.0	99.0	0.9

1. calc from Carb / protein / fat per each food.

$$10 \quad A = [[86, 0.0, 4.4, 68.0], \\ [1.2, 104.0, 52.0, 8.0], \\ [1.8, 135.0, 99.0, 0.9]]$$

$$\text{cal} = A \cdot \text{sum}(\text{axis}=0)$$

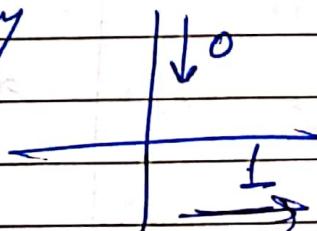
$$15 \quad [59, 239, 155.4, 76.9]$$

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

$$\text{per} = 100 \times \text{per} :$$

$$20 \quad [[9.9, 9.0, 2.88, 88.4], \\ [2.03, 43.51, 33.46, 10.40], \\ [3.05, 87.48, 63.70, 1.17]],$$

25 $A \cdot \text{sum}(\text{axis}=0)$ means python shade sum vertically



30 $\text{cal} \cdot \text{reshape}(1, 4)$ no need it's already is

L Broadcast Law

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

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \underbrace{\begin{bmatrix} 100 & 200 & 300 \end{bmatrix}}_{\text{1 row 3 columns}} \rightarrow \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$$

=
 $(m, n) \underset{\substack{* \\ |}}{+} (1, n)$ no copy n times (m, n)

and then multiplies the action (+, -, *, /)

$(m, n) \underset{\substack{* \\ |}}{+} (m, 1)$ no copy n times (m, n) & apply operation.

Python Rule of Thumbs

a = np.random.rand(5) ~~(X)~~ [rank 1 array]
 not use

a = np.random.rand(5, 1) ✓

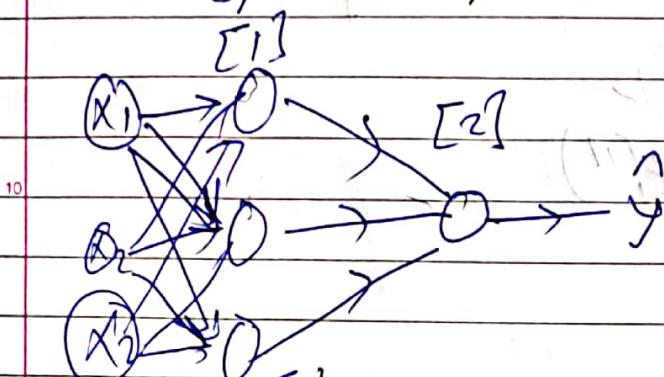
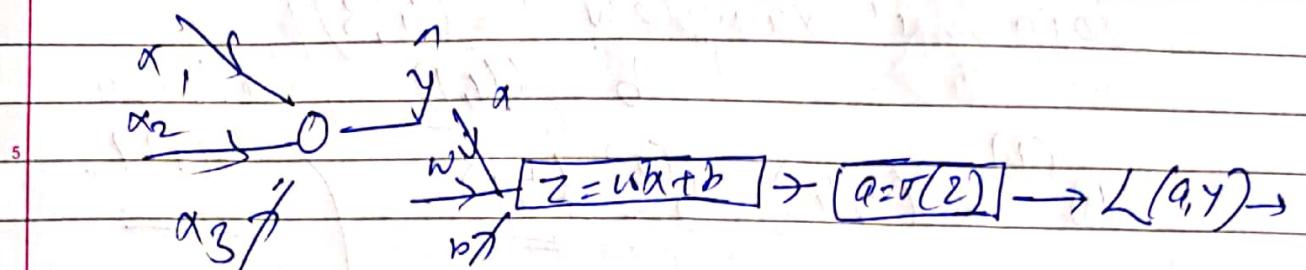
Use shape assert:

assert (a.shape == (5, 1))

Shallow Neural Network

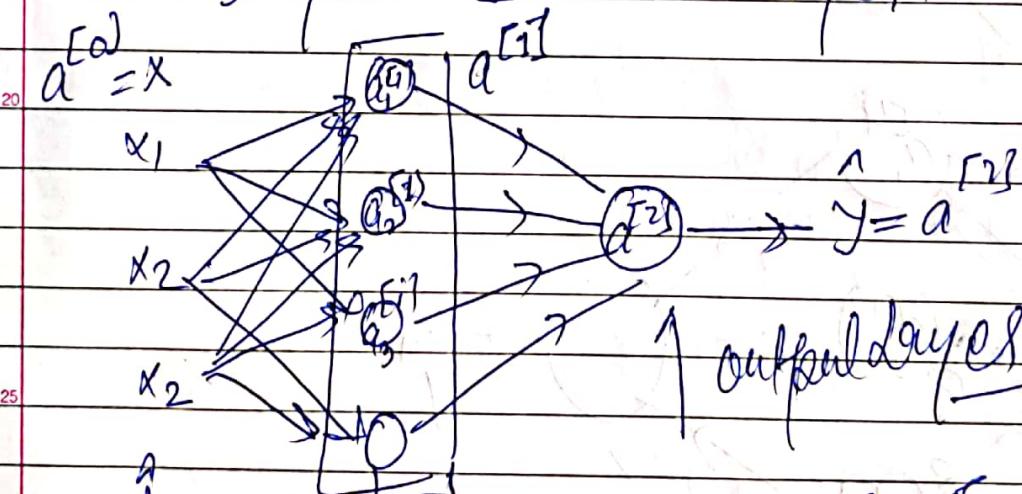
Camlin Page
Date / /

[Week - 03]



$w^{[1]}$ } allow us to find $E_1 = w^T x + b^{[1]}$
 $b^{[1]}$ } $a^{[1]} = \sigma(E_1)$

Superscript $[i]$ weight bias correspond to layer \Rightarrow i^{th} layer:



Input Layer Hidden Layer $a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \end{bmatrix}$

$a^{[i]}$ i^{th} layer activation
 " a layer, NN " / we don't count input layers

Hidden Layer Parameters $\vec{w}^{(1)}, \vec{b}^{(1)}$

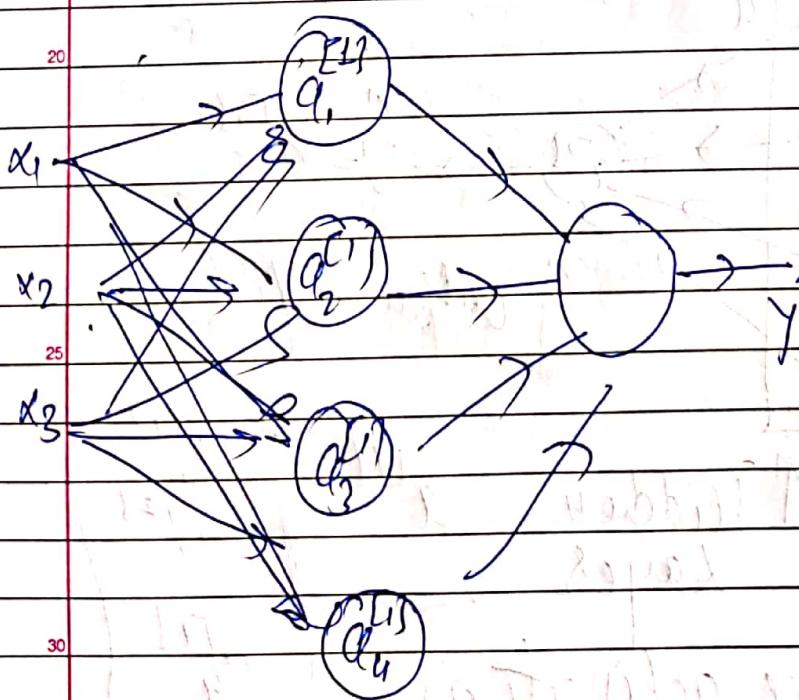
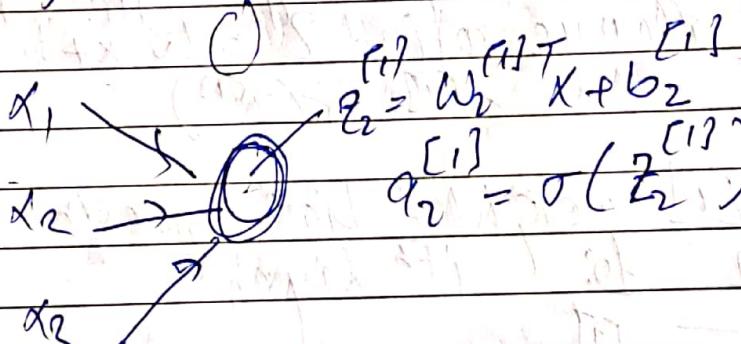
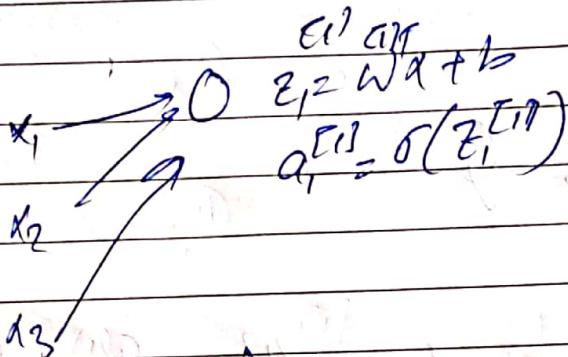
So in this example $w^{(1)} = (4, 3)$

$$b = (4, 1)$$

dimensions

$$w^{(2)} = (1, 4) \quad b^{(2)} = (1, 1)$$

dimensions



$$z_i^{[1]} = w_i^{[1]T} x + b_i \Rightarrow a_i^{[1]} = \sigma(z_i^{[1]})$$

each neuron ACT as one logistic unit
 $(w^T x + b)$

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \cdot \begin{pmatrix} w_1 \\ w_2 \\ w_3 \end{pmatrix} + b = (w, w_1, w_2) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

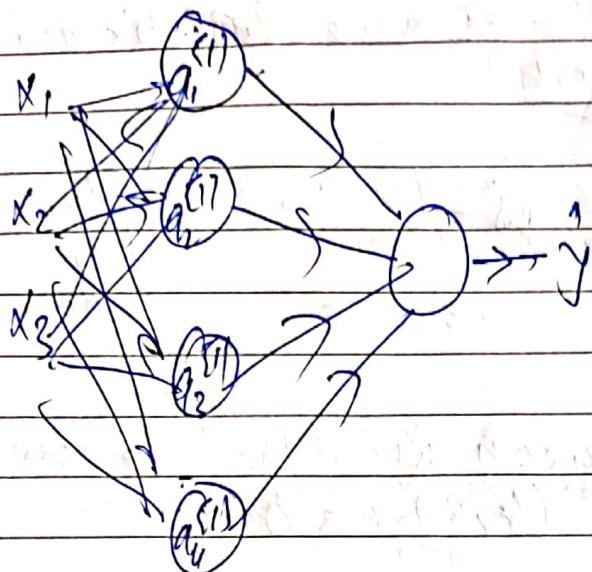
3×1

so for each neuron we'll put a row hence
 w will be $(4, 3)$ matrix.

$$z^{[1]} = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + \begin{pmatrix} b_{11}^{[1]} \\ b_{21}^{[1]} \\ b_{31}^{[1]} \\ b_{41}^{[1]} \end{pmatrix} = \underbrace{w^{[1]T} x + b^{[1]}}_{= -}$$

$$z = \begin{pmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{pmatrix} \Rightarrow a^{[1]} = \begin{pmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{pmatrix} = \sigma(z)$$

Given Input X



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

$$= w \cdot x + b$$

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

$$a^{(1)} = (4, 1) + (4, 1)$$

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

$$z^{(2)} = w^{(2)} a^{(1)} + b^{(2)} = (1 \times 4)(4 \times 1) + (1, 1) = (1 \times 1)$$

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

$$= (1 \times 1)$$

Across Multiple examples Vectorized

for $i = 1 \rightarrow m$

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

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

$$x = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix}$$

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

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

(n_x, m)

$m = \text{No. Train exm}$
 $n_x = \text{feature length.}$

Vectorizing It

$$w = \underbrace{\begin{pmatrix} w^{[1]} \\ w^{[2]} \end{pmatrix}}_{(4, 3)} \quad \underline{w^T = (3, 4)}$$

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

$$= b \cdot \underbrace{(x)}_{(3, 1)} +$$

$$\Rightarrow (4, 3)(3, m) + (4, 1)$$

$$z^{[1]} = (4, m) + (4, 1)$$

$$z^{[1]} = (4, m)$$

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

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

\leftarrow Across Train Ex.

$$a^{[1]} = \begin{bmatrix} a^{[1]}(1) & a^{[1]}(2) & \cdots & a^{[1]}(m) \end{bmatrix}$$

↑ different hidden unit
in a layer

$$W^{[1]} = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix} \quad X^{(1)} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$W^{[1]} X^{(1)} = \begin{bmatrix} \text{---} & 1 & 0 \\ \text{---} & 0 & 1 \\ \text{---} & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$= (4 \times 3) (3 \times 1) \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$W^{[1]} X^{(2)} = \begin{bmatrix} * \\ * \\ * \\ * \end{bmatrix} \quad W^{[1]} X^{(3)} = \begin{bmatrix} * \\ * \\ * \\ * \end{bmatrix}$$

$$X = \begin{bmatrix} X^{(1)} & X^{(2)} & X^{(3)} \\ 1 & 1 & 1 \end{bmatrix}$$

X

$$W^{[1]} X = W^{[1]} \begin{bmatrix} X^{(1)} & X^{(2)} & X^{(3)} \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}$$

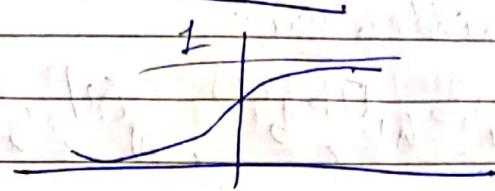
$$Z = \begin{bmatrix} Z^{E1} & Z^{E[1](2)} & Z^{E[1](3)} \\ 1 & 1 & 1 \end{bmatrix}$$

$$= Z$$

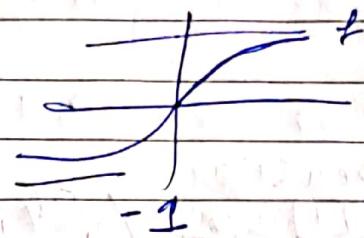
30

Activation functions

sigmoid (σ) =



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



tanh is almost superior than sigmoid fun

ReLU

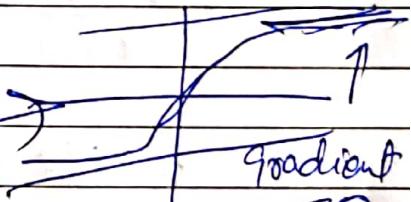
$$a = \max(0, z)$$

Derivatives of activation

$$\frac{d}{dz} g(z) = g(z)(1-g(z))$$

$g(z) = \text{sigmoid}$

$$\frac{d}{dz} g(z) = 1 - (\tanh(z))^2$$



$$g(z) = \text{ReLU}$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & z \geq 0 \end{cases}$$

In practice

$$\text{Random init: } w^{[1]} = np.random.randn(2, 2) * 0.01$$

Random init weights should be very small. tanh or sigmoid weight large gradient ≈ 0 .

Gradient Descent NN

1 hidden unit

Param: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

$h_x = h^{[0]}$
 $h^{[1]}$ hidden unit
 $h^{[2]} = 1$ output

dimensions

$$w^{[1]} = (h^{[1]}, h^{[0]})$$

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

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

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

Cost function, !:-

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

while (1) Σ

$$dw^{[1]} = \frac{\partial J}{\partial w^{[1]}}, db^{[1]} = \frac{\partial J}{\partial b^{[1]}}$$

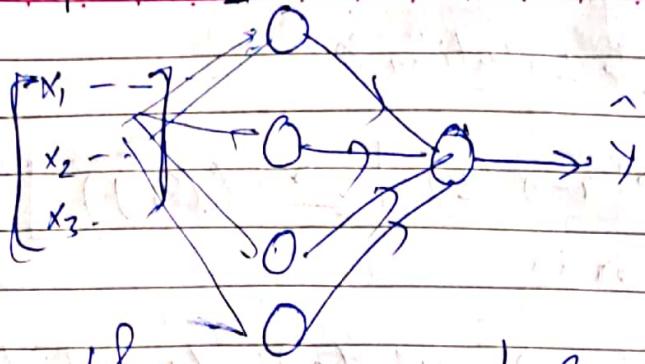
$$w^{[1]} = w^{[1]} - \alpha dw^{[1]}$$

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

$$w^{[2]} = w^{[2]} - \alpha dw^{[2]}$$

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

3



Ground Truth

forward prop

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

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

$$Z^{[1]} = (4, m) + (4, 1)$$

$$Z^{[1]} = (4, m)$$

$$\sigma(Z^{[1]}) = A^{[1]} = (4, m)$$

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

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

$$= (1, m) + (1, 1)$$

$$= (1, m)$$

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

~~$$A^{[2]} = [1, \dots, m]$$~~

Back propagation :-

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

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

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

keepdims=True)

Here it will give
(1, 4), (h^{[2]}, 1) matrix
so 1 * (1, 1) h^{[2]} = 1

$$dZ^{[1]} = W^{[2], T} dZ^{[2]},$$

~~$\underbrace{(h^{[1]}, m)}$~~
 ~~$\underbrace{dZ^{[2]}}$~~
 ~~$\underbrace{\# g^{[1]}(Z^{[1]})}$~~

Element wise product

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

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

Computing Gradients

(Backprop Intuition)

$$\begin{array}{c}
 \text{a} \\
 \overrightarrow{w} \\
 \downarrow b \\
 \boxed{z = w^T x + b} \rightarrow \boxed{a = \sigma(z)} \rightarrow \boxed{L(a, y)}
 \end{array}$$

d z d a d L

$$\begin{aligned}
 \frac{\partial L}{\partial a} &= \frac{\partial L}{\partial a} \quad L(a, y) = -y \log(a) - (1-y) \log(1-a) \\
 &= -\frac{y}{a} + \frac{1-y}{1-a} \\
 d\mathbf{z} &= a - y
 \end{aligned}$$

$$dz = da \cdot g'(z) =$$

$$\frac{\partial z}{\partial a} = \frac{\partial z}{\partial a} \left(\frac{\partial a}{\partial z} \right) \rightarrow \left(\frac{\partial g(z)}{\partial z} \right) = g'(z)$$

$$\begin{array}{c}
 \text{a}^{[1]} \\
 \overrightarrow{w^{[1]}} \\
 \downarrow b^{[1]} \\
 \boxed{z^{[1]} = w^{[1]} x + b^{[1]}} \rightarrow \boxed{a^{[1]} = \sigma(z^{[1]})} \rightarrow \boxed{z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}}
 \end{array}$$

d $z^{[1]}$ d $a^{[1]}$ d $w^{[1]}$ d $b^{[1]}$

$$\begin{array}{c}
 \text{a}^{[2]} = \sigma(z^{[2]}). \rightarrow \boxed{L(a^{[2]}, y)} \\
 d\mathbf{z}^{[2]} \qquad \qquad \qquad d\mathbf{a}^{[2]}
 \end{array}$$

$$dw^{[2]} = dt^{[2]} \cdot a^{[1]T}$$

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

$$w^{[2]} = (w^{[2]}, b^{[2]})$$

$$z^{[2]}_{id} = (z^{[2]}, 1).$$

$$z^{[1]} = (z^{[1]}, 1) = dt^{[1]}$$

$$\begin{aligned}
 dZ^{[1]} &= W^{[2]T} dZ^{[2]} * g^{(1)'}(z^{[0]}) \\
 &= (n^{(1)}, n^{(2)}) (n^{(2)} \times 1) * (n^{(1)}, 1) \\
 &= (n^{(1)}, 1) * (n^{(1)}, 1)
 \end{aligned}$$

5

Summary of Gradient Descent

$$\begin{aligned}
 dZ^{[2]} &= a^{[2]} - y \\
 dw^{[2]} &= dZ^{[2]} a^{[1]T} \\
 db^{[2]} &= dZ^{[2]}
 \end{aligned}$$

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{(1)'}(z^{[1]})$$

$$dw^{[1]} = dZ^{[1]} x^T \quad \downarrow \quad db^{[1]} = dZ^{[1]}$$

15

↳ Vectorized ↳

$$dZ^{[2]} = A^{[2]} - Y$$

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

20

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dZ^{[2]},$$

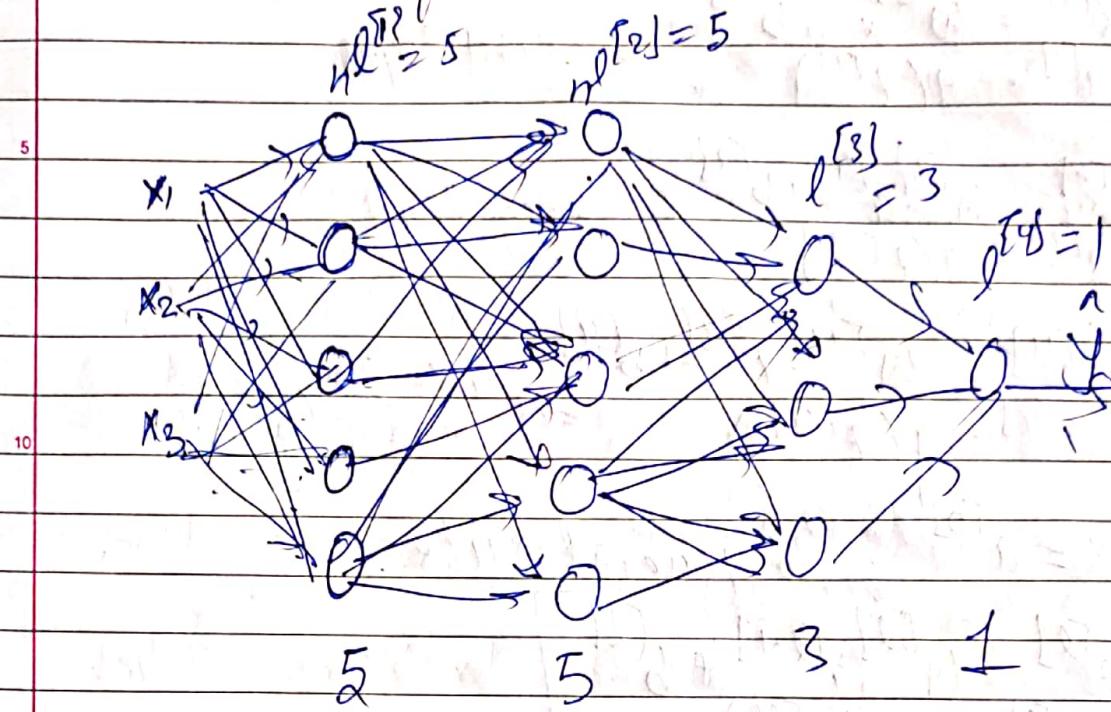
axis=1, keepdims=True)

$$\begin{aligned}
 dZ^{[1]} &= W^{[2]T} dZ^{[2]} \\
 &\ast g^{(1)'}(z^{[0]})
 \end{aligned}$$

30

Week - 04

Deep Neural Network



4 layer NN | $L = 4$ (No. of layers)
 $n^{[L]} = \# \text{ units in layer } (L)$

$$n^{[1]} = 5; n^{[2]} = 5; n^{[3]} = 3; n^{[4]} = n^4 = 1$$

$$h^{[0]} = h_x = 3$$

$a^{[l]}$ \Rightarrow activation in layer l
 $a^{[l]} = g^{[l]}(z^{[l]})$ | $w^{[l]}$ =

$x = \text{input feature} = a^{[0]}$
 $y = \text{output layer} = a^{[L]}$

Forward Propagation in a Deep Network

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

$$a^{[1]} = g(z^{[1]}) \quad = \text{first layer}$$

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

$$a^{[2]} = g(z^{[2]}) \quad = \text{second layer}$$

$$z^{[3]} = w^{[3]}a^{[2]} + b^{[3]}$$

$$a^{[3]} = g(z^{[3]}) \quad = \text{third layer}$$

$x = a^{[0]}$; So general Rule

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

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

Vectorize

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

$$x = A^{[0]}$$

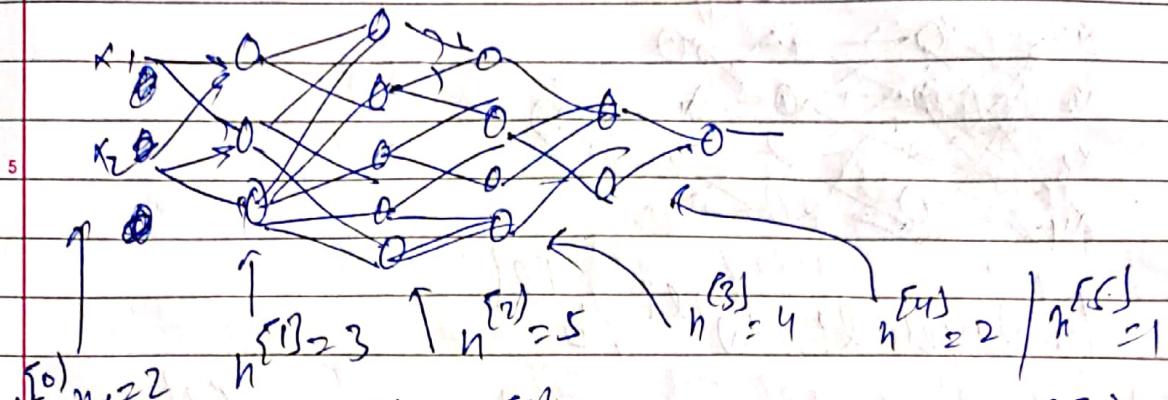
$$A^{[1]} = g(z^{[1]})$$

$$z^{[1]} = w^{[1]}A^{[0]} + b^{[1]}$$

25

30

Getting Dimension of Matrix Right



$$z^{(1)} = w^{(1)}x + b^{(1)}$$

$$(3,1) = (3,2)(2,1)$$

$$(h^{(1)}, 1) = (h^{(1)}, h^{(0)}) (h^{(0)}, 1)$$

$$w^{(1)} = (n^{(1)}, n^{(0)})$$

$$w^{(2)} = (n^{(2)}, n^{(1)}) (5,3)$$

$$b^{(2)} = b^{(1)}$$

$$(5,1) \quad (5,3)^T (3,1)$$

$$w^{(3)} = (4,5) \quad | \quad w^{(4)} = (2,4) \quad | \quad w^{(5)} = (1,2)$$

$$w^{(l)} = (n^{(l)}, n^{(l-1)})$$

$$b^{(l)} = (n^{(l)}, 1)$$

$$dw^{(l)} = (n^{(l)}, n^{(l-1)})$$

$$z^{(l)}, a^{(l)} \quad (n^{(l)}, m)$$

(Vectorize 18 lines) $dz^{(l)}, da^{(l)} : (n^{(l)}, m)$

$$z^{(1)} = w^{(1)}x + b^{(1)}$$

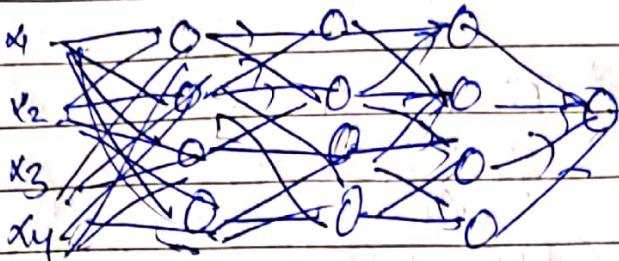
$$(h^{(1)}, 1) = (h^{(1)}, h^{(0)}) \rightarrow (h^{(0)}, 1).$$

$$z^{(1)} = w^{(1)}x + b^{(1)}$$

$$(h^{(1)}, m) = (h^{(1)}, n^{(0)}) (h^{(0)}, m) + (h^{(1)}, 1) \quad (\text{broadcasting})$$

Neural Network

(Building Blocks of Processing)

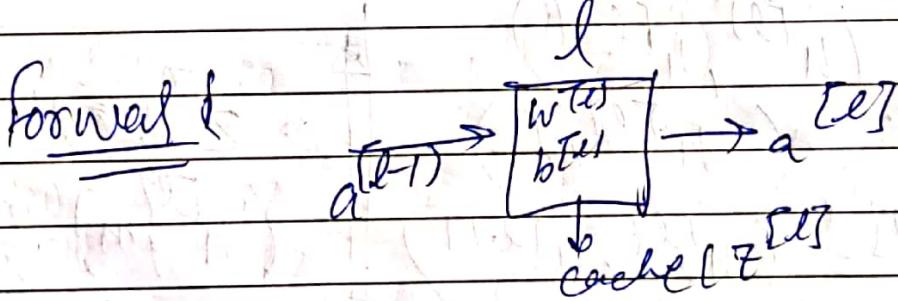


for layer l : $w^{[l]}, b^{[l]}$

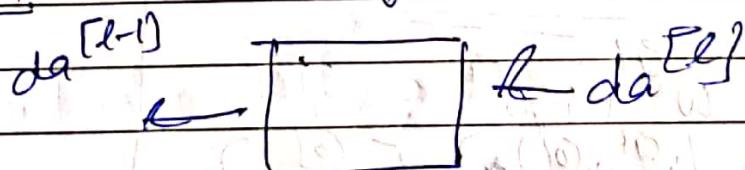
Forward: input $a^{[l-1]}$, output $a^{[l]}$

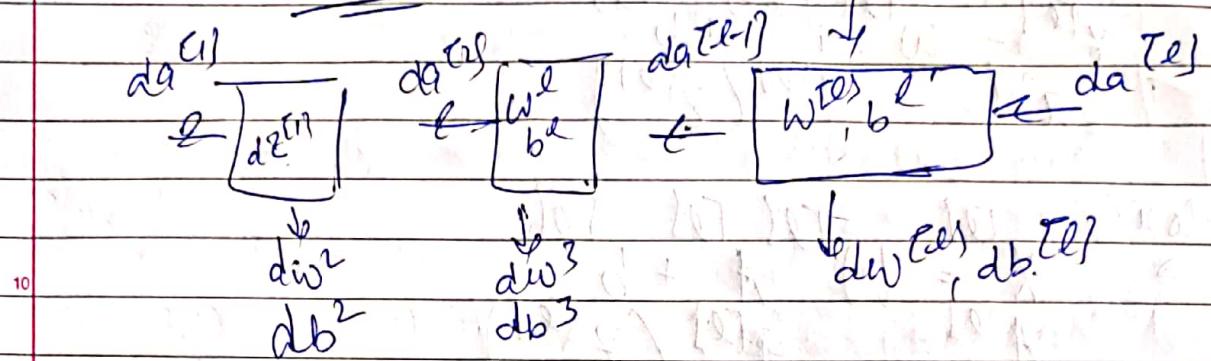
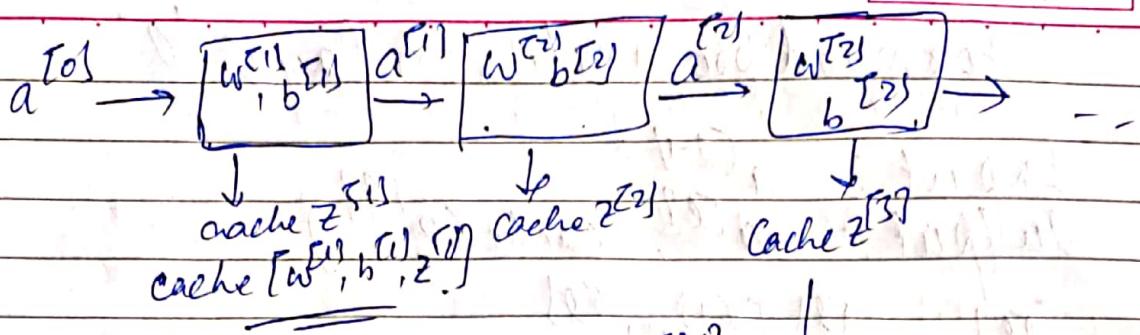
$$\begin{aligned} z^{[l]} &= w^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

Backward :- input $da^{[l]}$ output $da^{[l-1]}$
cache $z^{[l]}$



Backward





15

$$w^{T2} = w - \alpha dw^{T2}$$

$$b^{T2} = b^{T1} - \alpha db^{T2}$$

20

.....

25

.....

30

.....

f. To

Forward propagation for layer "l"

Input: $a^{[l-1]}$
 Output: $a^{[l]}$, cache($z^{[l]}$)

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

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

vectorize $a^{[l]}$ $z^{[l]}$ $A^{[l]}$ $S^{[l]}$

$$z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

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

Backward propagation for layer "l"

Input: $da^{[l]}$; Output: $da^{[l-1]}$, $dW^{[l]}$, $db^{[l]}$

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

$$dW^{[l]} = dz^{[l]} \cdot a^{[l-1]}$$

vector size

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

$$da^{[l-1]} = \cancel{dz^{[l-1]}} \cdot W^{[l]T} dz^{[l]}$$

$$dz^{[l-1]} = da^{[l-1]} + g^{[l-1]}'(z^{[l-1]})$$

$$dW^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} \text{np.sum}(dz^{[l]}, \text{axis=1, keepdim=True})$$

$$dA^{[l-1]} = W^{[l-1]} \cdot dZ^{[l]}$$

$$da^{[l]} = -\frac{y}{a} + \frac{(1-y)}{(1-a)}$$

$$dA^{[l]} = \left(-\frac{y^{(1)}}{a^{(1)}} + \frac{1-y^{(1)}}{1-a^{(1)}} \right) \cdot \dots \cdot \left(-\frac{y^{(m)}}{a^{(m)}} + \frac{1-y^{(m)}}{1-a^{(m)}} \right)$$

P.T.O

Z forward & Backward Prop

$$Z^{[1]} = W^{[1]} X + b^{[1]} \quad | \quad \text{d Backward Prop}$$

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

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

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

$$\begin{aligned} dZ^{[L]} &= A^{[L]} - Y^{[L]} \\ dW^{[L]} &= \frac{1}{m} dZ^{[L]} A^{[L]T} \end{aligned}$$

$$db^{[L]} = \frac{1}{m} npsum(dZ^{[L]}, axis=1, keepdims=True)$$

$$dZ^{[L-1]} = dW^{[L]T} dZ^{[L]} * g'^{[L]}(Z^{[L-1]})$$

$$dZ^{[1]} = dW^{[2]T} dZ^{[2]} * g'^{[2]}(Z^{[1]})$$

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

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

(axis=1, keepdims=True)