

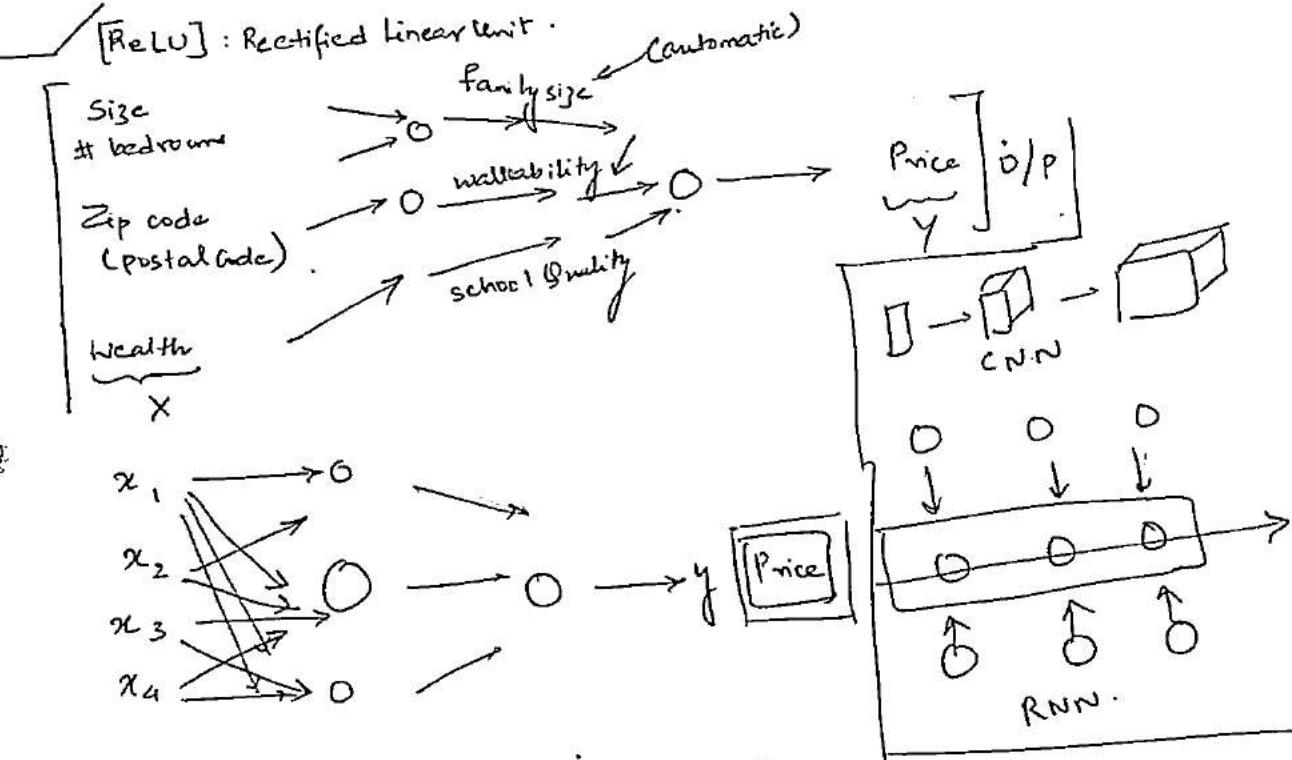
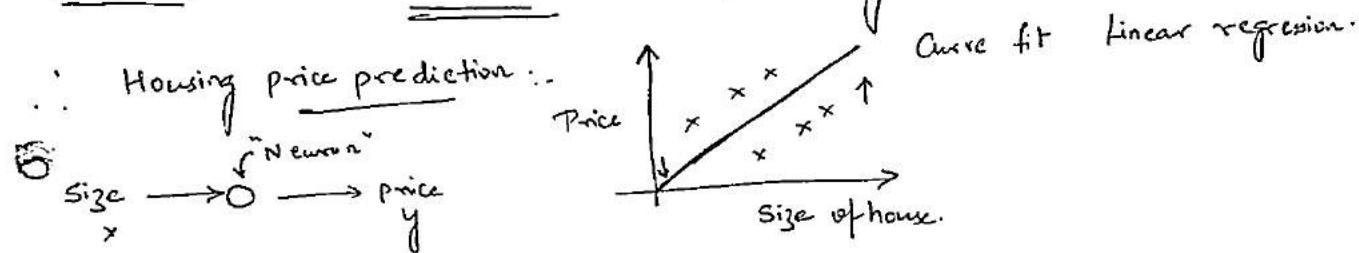
29/10/20
Target

Deep learning . A I

Coursera Specialization

- * Courses:-
1. Neural Networks and Deep Learning
2. Im
3.
4.
5.

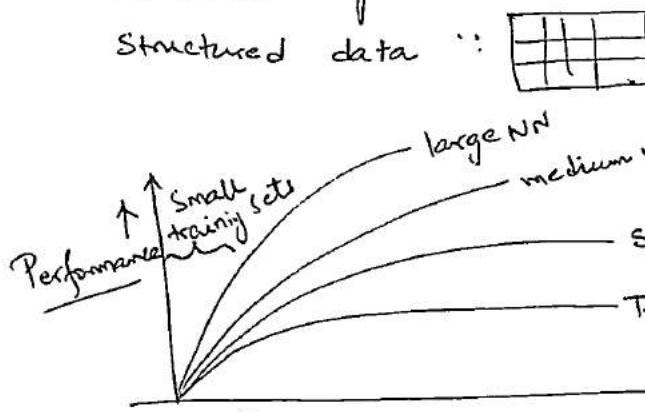
* Course 1: Neural Networks & Deep learning.



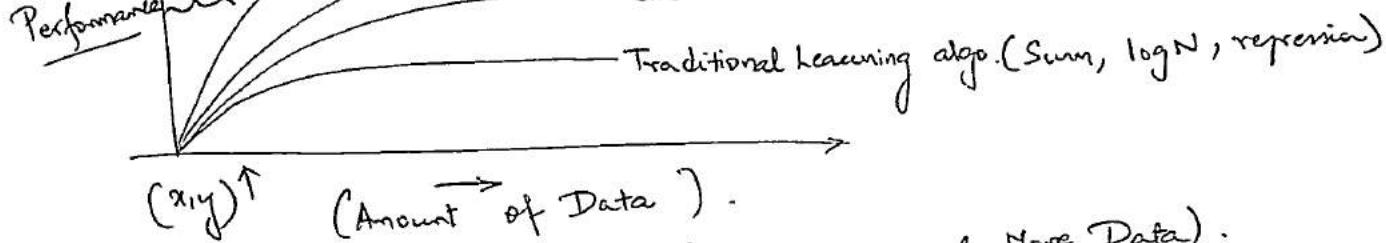
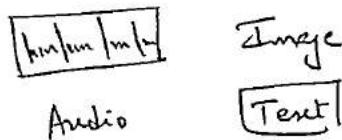
II Supervised learning		I/p(x)	O/p(y)	Application
Homefeat..		Price		Real Estate }
Ad, Userinfo		Click		Online Adver }
Image	Obj (...)			Photo tagging }
Audio	Text transcript			CNN (Convolutional)
Image/A-				Speech Recognition & Audi RNN (Recurrent)
				PTC-translations }
				Autonomous Driving }
				Custom/Hybrid

Supervised learning:

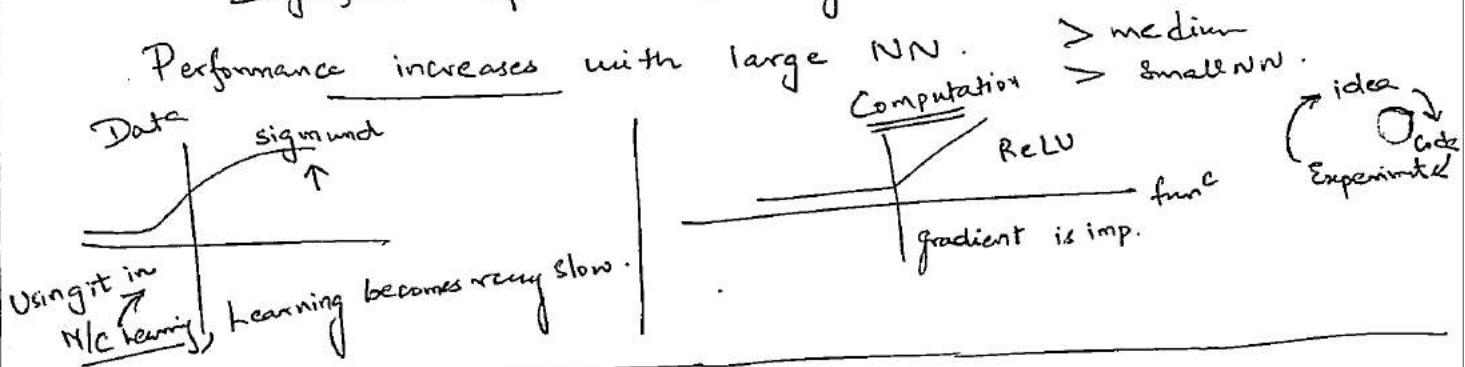
- 2 -



Unstructured Data



Digitization of Data (Collecting More & More Data).



INTRODUCTION TO NEURAL NETWORKS.

1. Neural Networks & Deep Learning -

Week 1 : Intro.

Week 2 : Basics of NN programming

Week 3 : One hidden layer NN

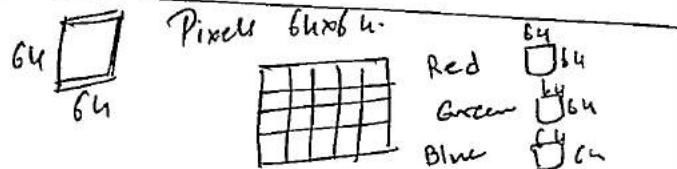
Week 4 : Deep Neural Networks.

Quiz :

1. AI is the new electricity :- Similar to electricity starting about 100 years ago, AI is transforming multiple industries.
2. Deep learning requires multiple iterations ; train models, experience is 2nd order.
3. $\frac{\partial}{\partial z}$ Increasing the size of NN, training set does not hurt algorithm performance & it may help significantly.

* Binary Classification :-

FPP/BPP : Forward/ Backward propagation



$$x = [:] = 64 \times 64 \times 3 = 12288 \quad \eta = \eta_x = 12288$$

$(x, y) \quad x \in \mathbb{R}^n, y \in \{0, 1\} \quad m$ training examples : $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$$x = [\begin{matrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(n)} \end{matrix}] \quad \underbrace{[\dots]}_{n} \quad \underbrace{[= x^{(m)}]}_{n}$$

$x \in \mathbb{R}^{n \times n}$

$$X.\text{shape} = (n, m) ; Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

Logistic Regression

$y \in \mathbb{R}^{l \times m}$

$y.\text{shape} = (l, m)$

Given X , want

$$\hat{y} = P(y=1|x)$$

$x \in \mathbb{R}^{n \times 1}$

(x dim vector).

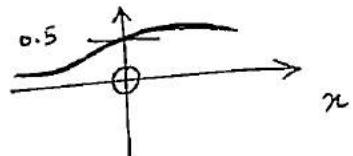
Par. $\boxed{w} \in \mathbb{R}^{n \times 1}$

$\boxed{b} \in \mathbb{R}$.

Output

$$\boxed{\hat{y} = w^T x + b}$$

$$0 \leq \hat{y} \leq 1$$



Sigmoid func.: $\hat{y} = \sigma(w^T x + b)$

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

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

if z large negative num, $\sigma(z) = \frac{1}{1+e^{-z}} \approx$

$$\approx \frac{1}{1+\text{Big num}} \approx 0.$$

alternative notation:-

$$x_0 = 1, x \in \mathbb{R}^{n \times 1}, \hat{y} = \sigma(\theta^T x)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_R \end{bmatrix} \begin{matrix} \left. \right\} b \leftarrow \\ \left. \right\} w \leftarrow \end{matrix}$$

Summary Parameters : $w \in \mathbb{R}^n, b \in \mathbb{R}, \hat{y} = w^T x + b, \hat{y} = \sigma(w^T x + b)$

$$\boxed{\hat{y} = \sigma(\theta^T x); \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_R \end{bmatrix} \begin{matrix} \left. \right\} b \leftarrow \\ \left. \right\} w \leftarrow \end{matrix}}$$

Logistic Regression Cost function

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

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

* loss (error) func : $L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

(efficiency of Alg.)

how good \hat{y} is ?? given by

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

* $L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$ if $y=0: L(\hat{y}, y) = -\log (1-\hat{y})$

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

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

Summary:

∴ w, b minimize cost func;

loss func for efficiency of algorithm

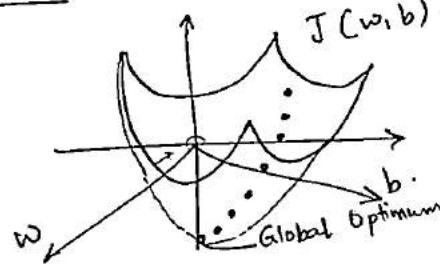
Logistic
Regression

Note: The loss func computes the error for a single training example, the cost func is the average of the loss funcs of the entire training set.

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

Gradient Descent: - (train / learn w, b on training set).

until now:



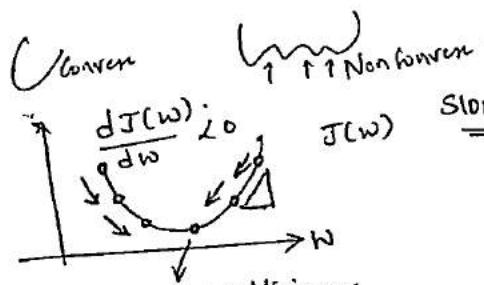
Cost func: $J(w, b)$

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

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

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

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



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

learning rate
update "dw".

$$w := w - \alpha \partial w \quad \text{Gradient Descent Update.}$$

$J(w, b)$
Update rule
for Gradient
Descent

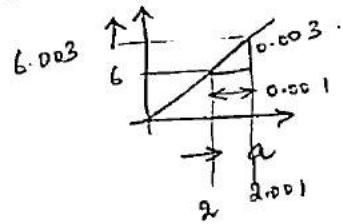
$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

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

∂ = partial derivative

∴ A convex func has single local optima

Derivatives :-



$$f(a) = 3a$$

$$a = 2, f(a) = 6$$

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

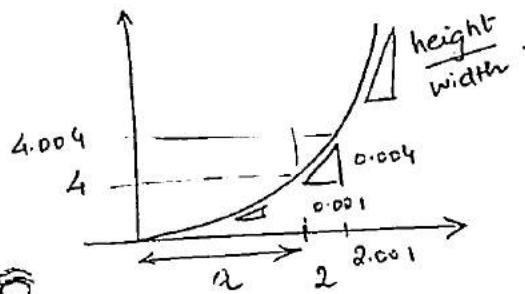
Slope derivative of $f(a)$.
at $a = 2$ is 3.

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

Slope at $a = 5$ is also 3

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

* Examples : Intuition about derivatives :



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

$$\frac{d}{da} f(a) = 4, \text{ when } a = 2$$

\therefore if $a = 5$ $f(a) = 5$
 $a = 5.001$ $f(a) = 25.010.$

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

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

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

$$\text{as } a \uparrow \quad \frac{d}{da} f(a) = \uparrow \uparrow \uparrow$$

$$a = 2 \quad f(a) = 4$$

$$a = 2.001 \quad f(a) \approx 4.004$$

$$f(a) = a^3 \quad \frac{d}{da} f(a) = 3a^2; 3 \times a^2 \quad a = 2 \quad f(a) = 8$$

$$\rightarrow 3 \times 2.001^2 \quad a = 2.001 \quad f(a) \approx 8.012.$$

$$f(a) = \log_e(a) \quad \ln(a)$$

$$\frac{d}{da} f(a) = \frac{1}{a} \quad a = 2 \quad f(a) = 0.69315$$

$$a = 2.001 \quad f(a) = 0.69365$$

$$\frac{d}{da} f(a) = \frac{1}{2}$$

Derivative is slope of a line;

* Slope is constant for linear lines whereas varies for exponential, log etc.

*

Computation graph :- FPP/BPP. $J(a, b, c) = 3(a + bc)$.

$$\therefore u = bc$$

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

$$\begin{array}{l} s = a \\ 3 = b \\ 2 = c \end{array}$$

$$\begin{array}{c} 6 \\ \rightarrow [u = bc] \rightarrow [v = a + u] \rightarrow [E = 3v] \\ \downarrow \qquad \downarrow \qquad \downarrow \\ 33 \end{array}$$

Note: One step of backward propagation on a computation graph yields derivative of final output variable.

Derivatives with a computation graph:-

* Derivation of $\frac{dL}{dz}$

8, 6, 10, 7, 9, 11, 6, 3, 7,

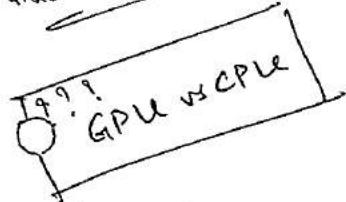
Time 12:00

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

* Vectorization:-

$$z = w^T \times b.$$

if vectorize your code
Deep learning runs
much faster



GPU

CPU } SIMD - single instruction multiple data.

Basics of Neural Network Programming

Deep learning, AI.

Normal

```
a = np.random.rand(1000000)
b = " "
tic = time.time()
for i in range(100000):
    c += a[i]*b[i]
toc = time.time()
```

4.81.31102039 ms | 1.5081 ms.

Vectorize

300 times faster than normal

```
a = np.random.rand(1000000)
b = " "
tic = time.time()
c = np.dot(a, b)
toc = time.time()
```

GPU (Graphical P.U. has many cores & used for parallel computing (Simult. many things))
CPU (few cores & less parallelization)

Note : if you use In-Built func, it enables python to efficient parallelism]



Examples :-

AVOID Explicit "FOR" Loops.

$$u = Av,$$

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

$$\text{Normal} \quad u = np.zeros(n, 1).$$

Vectorize

but

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

* Vectors & Matrix valued functions

Apply the exponential operation on every element of a matrix/vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

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

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} , 1/v.

$$\therefore dw = np.zeros((n, 1)) \quad (\text{dimensional vector})$$

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

$$dw = dw/m$$

$$|dw| = m$$

for training examples Vectorizing Logistic Regression

$$\rightarrow z^{(1)} = w^T x^{(1)} + b \quad z^{(2)} = w^T x^{(2)} + b \quad z^{(3)} = w^T x^{(3)} + b \\ a^{(1)} = \sigma(z^{(1)}) \quad a^{(2)} = \sigma(z^{(2)}) \quad a^{(3)} = \sigma(z^{(3)}) \\ \therefore x = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_m^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_m^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_m^{(m)} \end{bmatrix} \in \mathbb{R}^{n_x \times n_m}$$

$(z = np.dot(w.T, x) + b)$ ("Broadcasting" in Python)

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

VECTORIZING LOGISTIC REGRESSION GRADIENT OUTPUT

$$\frac{\partial z^{(1)}}{\partial w} = a^{(1)} - y^{(1)} \quad \frac{\partial z^{(2)}}{\partial w} = a^{(2)} - y^{(2)} \dots \\ \frac{\partial z}{\partial w} = [d z^{(1)} \ d z^{(2)} \ \dots \ d z^{(m)}]. \quad A = [a^{(1)} \dots a^{(m)}] \cdot Y = [y^{(1)} \dots y^{(m)}]$$

$\frac{\partial z}{\partial b} = A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots].$

$$\frac{\partial b}{\partial b} = \frac{1}{m} \sum_{i=1}^m \frac{\partial z^{(i)}}{\partial b} \\ = \frac{1}{m} \text{np.sum}(dz).$$

$\frac{\partial w}{\partial w} = \frac{1}{m} X \cdot dz^T.$

$$= \frac{1}{m} \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_m^{(1)} \\ 1 & x_1^{(2)} & \dots & x_m^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & \dots & x_m^{(m)} \end{bmatrix} \begin{bmatrix} d z^{(1)} \\ d z^{(2)} \\ \vdots \\ d z^{(m)} \end{bmatrix} \\ = \frac{1}{m} \begin{bmatrix} x_1^{(1)} d z^{(1)} + \dots + x_m^{(1)} d z^{(m)} \\ x_1^{(2)} d z^{(1)} + \dots + x_m^{(2)} d z^{(m)} \\ \vdots \\ x_1^{(m)} d z^{(1)} + \dots + x_m^{(m)} d z^{(m)} \end{bmatrix} \in \mathbb{R}^{n_x \times 1}$$

$$\frac{\partial w}{\partial w} = 0 \quad \frac{\partial b}{\partial b} = 0 \\ \frac{\partial w}{\partial w} + x^{(1)} \frac{\partial z^{(1)}}{\partial w} \quad \frac{\partial b}{\partial b} + d z^{(1)} \\ \frac{\partial w}{\partial w} + x^{(2)} \frac{\partial z^{(2)}}{\partial w} \quad \frac{\partial b}{\partial b} + d z^{(2)} \\ \vdots \quad \vdots \\ \frac{\partial w}{\partial w} + d z^{(m)} \quad \frac{\partial b}{\partial b} + d z^{(m)}$$

Broadcasting in Python (faster).

$$A = np.array([[[\dots, \dots, \dots], [\dots, \dots, \dots], [\dots, \dots, \dots]]]) \\ \text{cal} = A.sum(axis=0) \\ \text{percentage} = 100 * A / \text{cal.reshape}(1, 1, 1)$$

Matlab/Octave: bsxfun

$$= \begin{bmatrix} \dots & \dots & \dots & \dots \\ | & | & | & | \\ \dots & \dots & \dots & \dots \end{bmatrix}$$

$\downarrow^0 \text{axis}$
 $\rightarrow \text{axis } 1$

$\text{percentage} = 100 * A / (\text{cal.reshape}(1, 1, 1))$

Ex. ①

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 202 \\ 303 \\ 404 \end{bmatrix}.$$

GENERAL PRINCIPLE of BoC

$$(m, n) \xrightarrow{*} (1, n) \rightsquigarrow (m, n)$$

$$(m, 1) \xrightarrow{*} (R \rightsquigarrow m, 1)$$

② $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 1,000 & 2,000 & 3,000 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 204 & 205 & 206 \end{bmatrix}. \quad (m, 1)$

bnt Python realisiert (m, n)

$$(1, m) \xrightarrow{*} (R \rightsquigarrow 1, m)$$

A Note on python/numpy vectors.

never use data structures

$a = np.random.random(5)$ \times $\xrightarrow{\text{rank 1 array}}$

$a = np.random.random(5, 1)$. $\xrightarrow{\text{rank 1 array}}$

$\xrightarrow{\text{point } (a.T) \text{ a transpose.}}$

$\xrightarrow{\text{np.dot } (a, a.T) \rightarrow \text{outerproduct of a vector}}$

$\xrightarrow{\text{value}}$

$\xrightarrow{\text{use in NN}}$

Note Never use rank 1 array.
Always use vector

$a = np.random.random(5, 1) \rightarrow a.shape = (5, 1)$ column vector

$a = np.random.random(1, 5) \rightarrow a.shape = (1, 5)$ row vector

$\xrightarrow{\text{Ouchshape an array}}$

$a = a.reshape((5, 1))$

$\xrightarrow{\text{assert } (a.shape = (5, 1))}$ inexpensive assert

Tour of Jupyter

Explanation of logistic regression Cost func

* $y | x$ y given x . Note: Minimizing the loss corresponds with maximizing $\log P(y|x)$.

Exercises: Week 2

Q1: A neuron computes a linear func^c ($z = w^T x + b$) followed by an activation func^c
We generally say that the output of a Neuron is a $g(w^T x + b)$ where g is the activation func^c $\hat{y} = g(w^T x + b) ; g(z) = \frac{1}{1+e^{-z}}$
~~(Sigmoid, tanh, ReLU, ...).~~

Q2. Logistic Loss :- $L^{(c)}(\hat{y}^{(c)}, y^{(c)}) = -(y^{(c)} \log(\hat{y}^{(c)}) + (1 - \hat{y}^{(c)}) \log(1 - \hat{y}^{(c)}))$.

Q3: $\text{img} = (32, 32, 3)$ array $\approx x = \text{img.reshape}((32*32*3, 1))$, column vector

Q4 $a = np.random.random(2, 3); c = a+b; c.shape = (2, 3)$ [Broadcasting of b . to match a].

Q5 $a = np.random.random(4, 3); b = " " (3, 2); c = a * b;$ "Error" [The '*' operator in Numpy, indicates element-wise multiplication.]

If it is different from "np.dot()". If you would try "c = np.dot(a, b)" you would get $c.shape = (4, 2)$

Q6. Suppose you have n_x input features per example. $X = [x^{(1)} \ x^{(2)} \dots x^{(m)}]$.

What is the dimension of X ? (n_x, m)

???

Q7 $a = np.random.rand(12288, 150)$; $c = np.dot(a, b)$; $c.shape = (12288, 45)$.
 $b = \dots (150, 45)$.

Q8. $a.shape = (3, 4)$. for i in range (3):

$b.shape = (4, 1)$ for j in range (4)

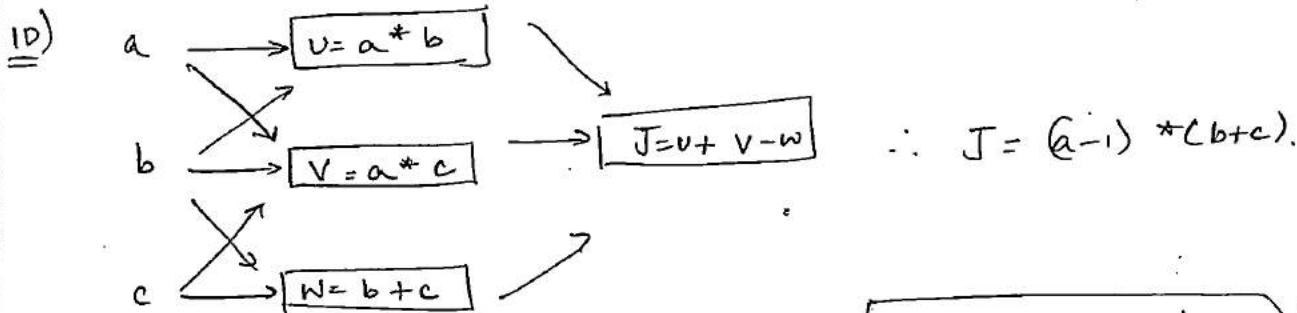
$$c[i][j] = a[i][j] + b[j].$$

Vectorize $C = a + b^T$.

Q9. $a = np.random.rand(3, 3)$. $c = a * b$

$b = np.random.rand(3, 1)$,

This will invoke broadcasting, so b is copied three times to become $(3, 3)$ and $*$ is an element-wise product so $c.shape$ will be $(3, 3)$.



* Programming exercise. Q. Build a func returns

math.exp(x)

$$\text{Sigmoid}(x) = \frac{1}{1+e^{-x}}$$

Logistic
func

Note: We rarely use "math" library in deep learning because the inputs of the func are real numbers. In deep learning we mostly use matrices and vectors.

This is why "numpy" is more useful.

Furthermore, if x is a vector, then Py operation such as $S = x + 3$ or $S = \frac{1}{x}$ will output S as a vector of the same size as x .

Note: ML/DL Normalize our data because gradient Descent converges faster after Normalization.

* gradient of Sigmoid func = slope of sigmoid func. $\frac{x}{\|x\|} \equiv np.linalg.norm(x, axis=1, keepdims=True)[S]$

Broadcasting changes matrix dimensions. $\int np.dot() \rightarrow \text{matrix-matrix multiplication}$

$mp.multiply() \rightarrow \text{element wise multiply}$

Programming Exercise Part 2

Centre & Standardize your
Data Set in Machine Learning

Logistic Regression to identify cats.

255 (maximum value of a pixel channel).

* m-train images labelled as Cat ($y=1$) & non-cat ($y=0$)

* $(\text{num_px}), (\text{num_px}, 3)$ 3 → 3 channels (RGB).

* Logistic regression is actually a very simple neural network

$$w^T x^{(i)} + b \mid 0$$

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

* The main steps for building a Neural Networks are:

1. Define the model structure (such as number of input features).

2. Initialize the models parameters.

3. loop: - Calculate current loss (forward propagation).

- calculate current gradient (backward propagation).

- update parameters (gradient descent).

$$W = np.zeros((dim, 1))$$

You often build 1-3 separately and integrate them into one func we call `modelC`.



Propagate():

W, b, X, Y .

$$A = \sigma(W^T X + b) = (a^{(1)}, a^{(2)}, \dots, a^{(m)})$$

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

$$\frac{\partial J}{\partial W} = \frac{1}{m} X(A-Y)^T$$

$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)})$$

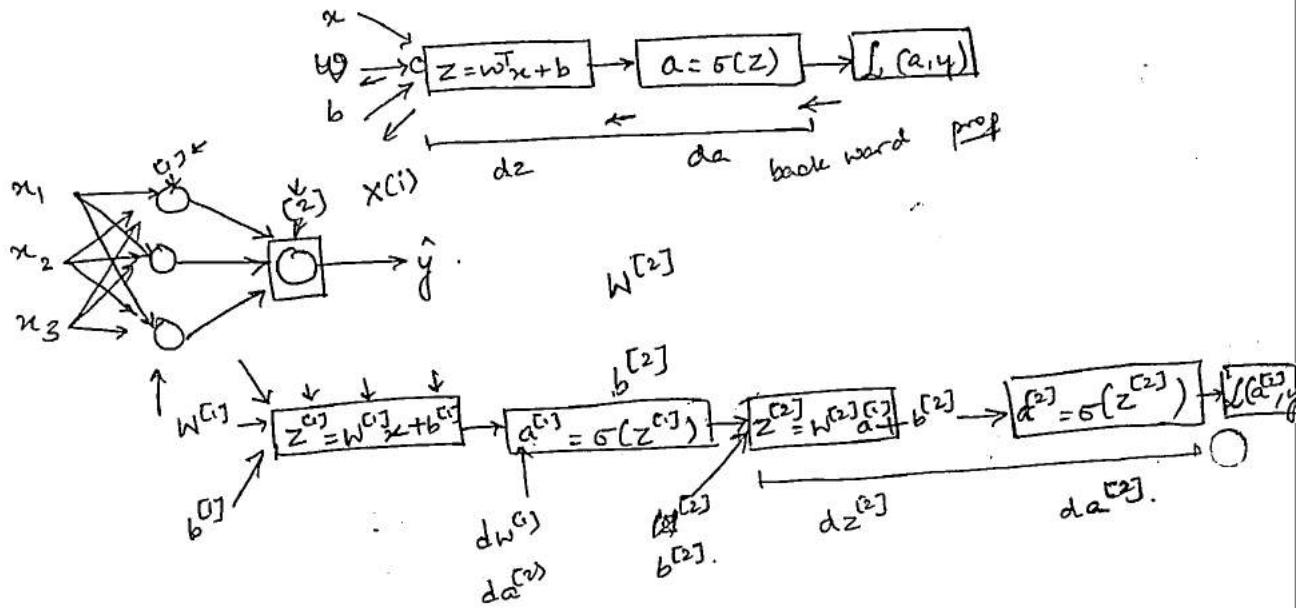


$$\theta = \theta - \alpha \delta \theta$$

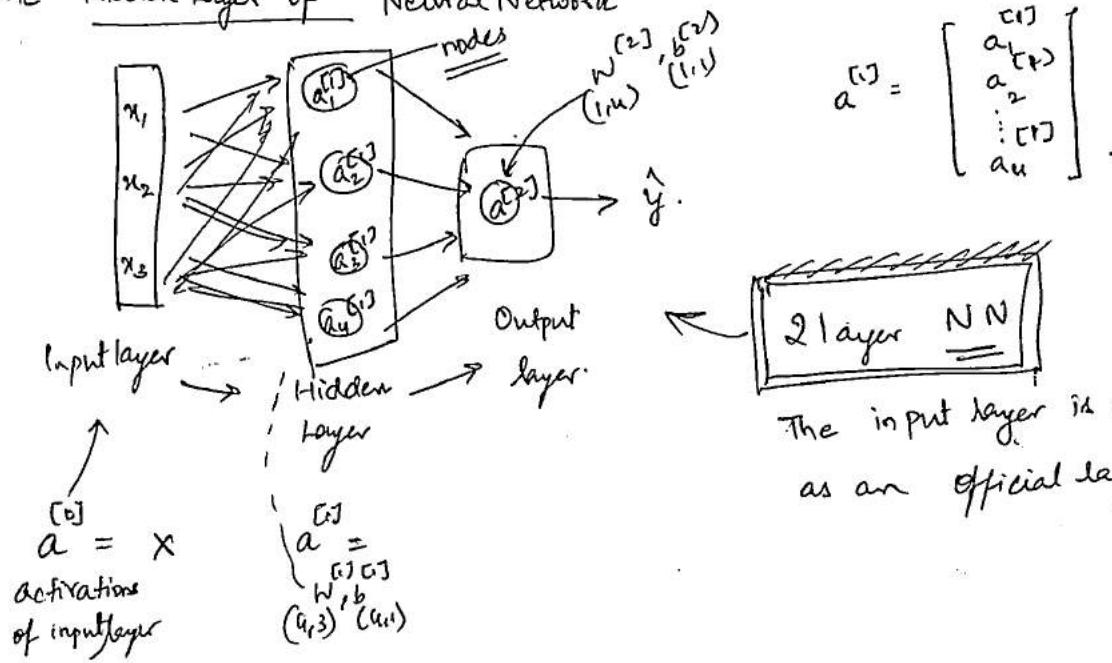
α = learning rate

~~Glimpse~~

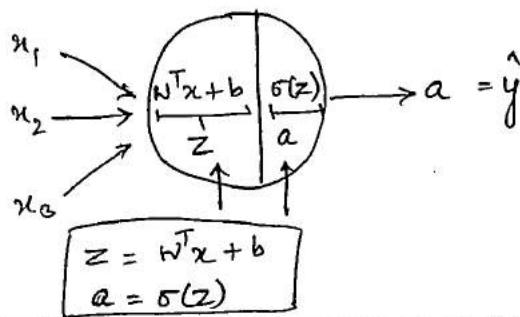
$$x_1 \rightarrow O \rightarrow \hat{y} = \varphi.$$



One Hidden layer of Neural Network

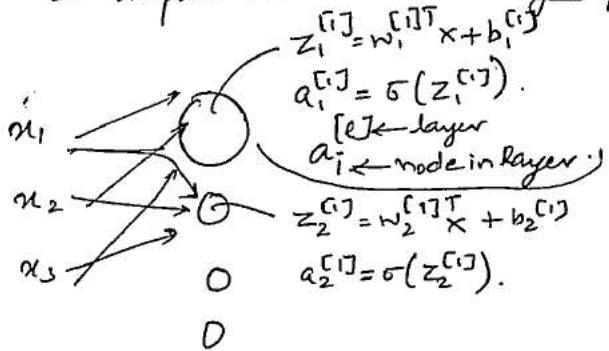


Computing a NN output : - 2 steps



1. Compute z

2. Compute Activation(a) from sigmoid/ σ



$$x_1 \rightarrow z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$$

$$w_2 \rightarrow a_2^{[1]} = \sigma(z_2^{[1]}).$$

* Neural Network Representation

$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, a_1^{[1]} = \sigma(z_1^{[1]}). \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, a_2^{[1]} = \sigma(z_2^{[1]}). \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, a_3^{[1]} = \sigma(z_3^{[1]}). \\ z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, a_4^{[1]} = \sigma(z_4^{[1]}). \end{aligned}$$

$$z^{[1]} = \underbrace{\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix}}_{(4,3)} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_{(3,1)} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}}_{(4,1)} = \underbrace{\begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ w_3^{[1]T} x + b_3^{[1]} \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix}}_{(4,1)} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

Given input x :

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]}).$$

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

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

$$\rightarrow z^{[1]} = w^{[1]T} x + b^{[1]}.$$

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

$$\rightarrow z^{[2]} = w^{[2]T} a^{[1]} + b^{[2]}.$$

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

9, 10, 11, 12, 13, 14

* Vectorize across multiple examples :-

$$\begin{aligned} x &\longrightarrow a^{[2]} = \hat{y}. \\ x^{[1]} &\longrightarrow a^{[2](1)} = \hat{y}^{(1)}. \\ x^{[2]} &\longrightarrow a^{2} = \hat{y}^{(2)}. \\ \vdots &\vdots \\ a^{[2](i)} & \text{example } i. \\ & \text{layer 2.} \end{aligned}$$

$$X = \begin{bmatrix} 1 & x^{[1]} & x^{[2]} & \dots & x^{[m]} \end{bmatrix} \quad Z^{[1]} = \begin{bmatrix} 1 & z^{1} & z^{[1](2)} & \dots & z^{[1](n)} \end{bmatrix}$$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$
 $(n, m) \quad \text{training examples} \quad \text{hidden units.}$

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

$$\left\{ \begin{array}{l} z^{[1]} = w^{[1]} x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}). \end{array} \right.$$

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

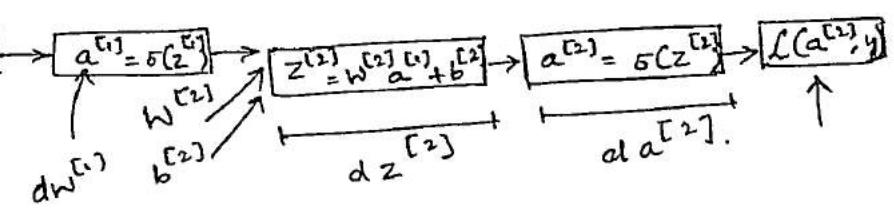
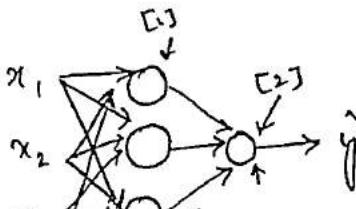
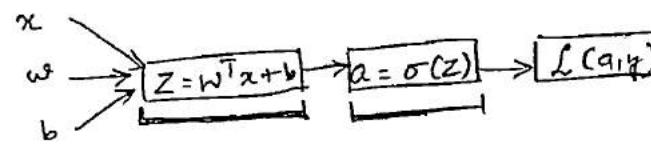
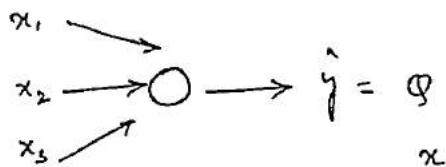
$$\left\{ \begin{array}{l} z^{[1](i)} = w^{[1]} x + b^{[1]} \\ a^{[1](i)} = \sigma(z^{[1](i)}) \\ z^{[2](i)} = w^{[2]} a^{[1](i)} + b^{[2]} \\ a^{[2](i)} = \sigma(z^{[2](i)}). \end{array} \right. \quad \square$$

$$\begin{aligned} z^{[1]} &= w^{[1]} x + b^{[1]}. \\ A^{[1]} &= \sigma(z^{[1]}). \end{aligned}$$

* Explanation for Vectorized implementation

$$\therefore \text{first training example} \quad z^{1} = w^{[1]} \underline{x^{(1)}} + \underline{b^{[1]}} \quad z^{[1](2)} = w^{[1]} \underline{x^{(2)}} + \underline{b^{[1]}}$$

$$w^{[1]} = \begin{bmatrix} \dots \\ \dots \end{bmatrix} \quad \square$$



* Activation (Network) functions:-

Better
funk \rightarrow Sigmoid funk

$$z^{[1]} = w^{[1]} x + b^{[1]} \quad \text{until now Activation func}$$

$$a^{[1]} = \sigma(z^{[1]}) \rightarrow g(z^{[1]}) \quad \text{is the sigmoid func}$$

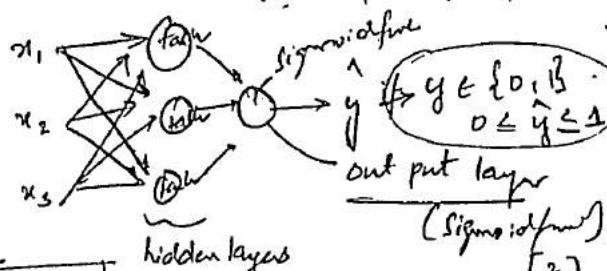
$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}.$$

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

Nonlinear funk (not sigmoid).

Now Activation func is tanh func

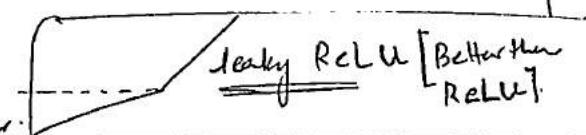
$$g(z^{[1]}) = \tanh(z^{[1]}).$$



Note

hidden layers have data closer to zero mean so, tanh better for Output layer (sigmoid) better.

[1].



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

ReLU

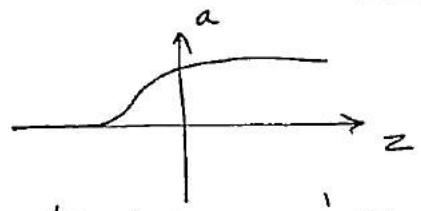
Rectified Linear Unit

(Default choice of Activation

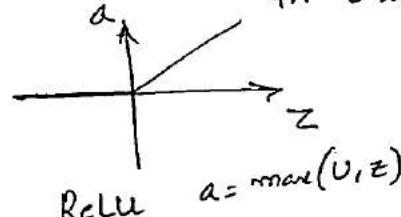
func for hidden layers if not sure then ... tanh or ReLU

- NN is much faster when using ReLU Activation func than tanh or sigmoid Activation func

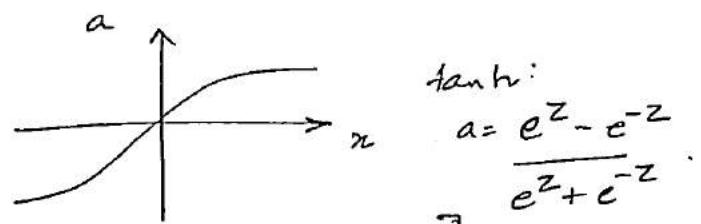
* Pros & Cons of Activation func. - 2 - 23, 24, 25, 26, 27, 28, 29, 30, 1, 2, 3, 4, 5, 6, 7.
15 days.



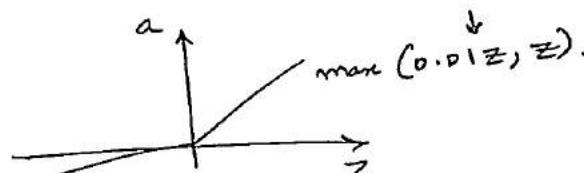
Sigmoid : $a = \frac{1}{1 + e^{-z}}$
 [Never for hidden layers & always for Output layers].



ReLU $a = \max(0, z)$
 (Default) for hidden layers.



[best for hidden layers]



leaky ReLU $a = \max(0.01z, z)$
 [sometimes better than ReLU].

* Why do we need a N.L Activation func. - ?

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

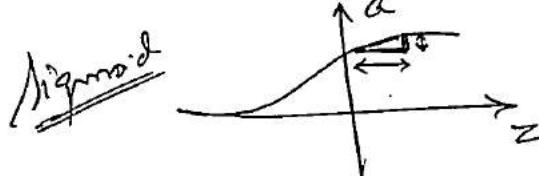
$$a^{[i]} = g^{[i]}(z^{[i]})$$

" $g(z) = z$
 linear orientation function"

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

∴ A linear Hidden layer
 is more or less
 useless. N.L is better

* Derivatives of Activation Functions



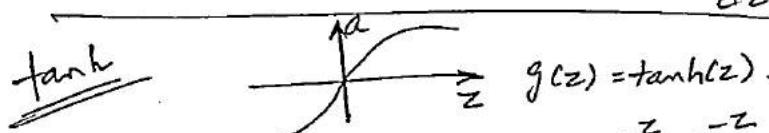
$$g(z) = \frac{1}{1 + e^{-z}}$$

$\therefore \frac{d}{dz}(g(z)) = \text{slope of } g(z) \text{ at } z$

$$= \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}} \right).$$

$$= g(z)(1 - g(z)).$$

$$= a(1 - a).$$



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

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

$$= 1 - (\tanh(z))^2$$

$$z = 10 \quad \tanh(z) \approx 1$$

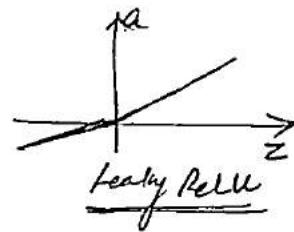
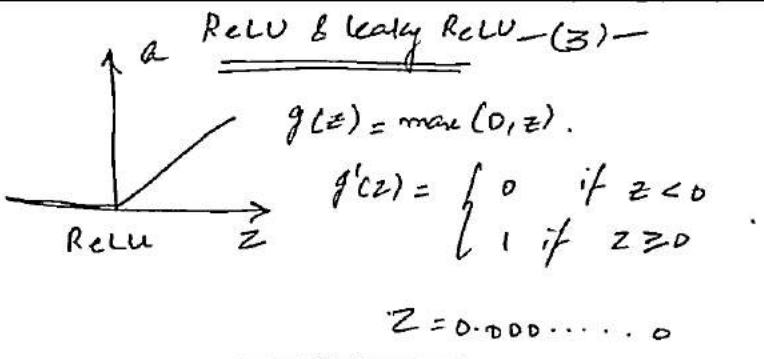
$$g'(z) \approx 0.$$

$$z = -10 \quad \tanh(z) \approx -1$$

$$g'(z) \approx 0.$$

$$z = 0 \quad \tanh(z) = 0$$

$$g'(z) = 1.$$



$$g(z) = \max(0.01z, z)$$

Gradient Descent for Neural Networks:-

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Parameters :- $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$.
 $(n^{[0]}, n^{[1]}) (n^{[1]}, 1) (n^{[2]}, n^{[1]}) (n^{[2]}, 1)$.

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

○ Gradient Descent :

Repeat { Compute predicts $(\hat{y}^{(i)}, i = 1, \dots, m)$.

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

Update rule $w^{[1]} := w^{[1]} - \alpha dw^{[1]}$, $b^{[1]} := b^{[1]} - \alpha db^{[1]}$, Update Rule for Gradient Descent

○ Formulas for computing derivatives:

Forward Propagation:

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

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

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

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

Correct vectorised Implementation
of forward propagation for layer i ; $1 \leq i \leq L$:

$$z^{[i]} = w^{[i]}A^{[i-1]} + b^{[i]}.$$

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

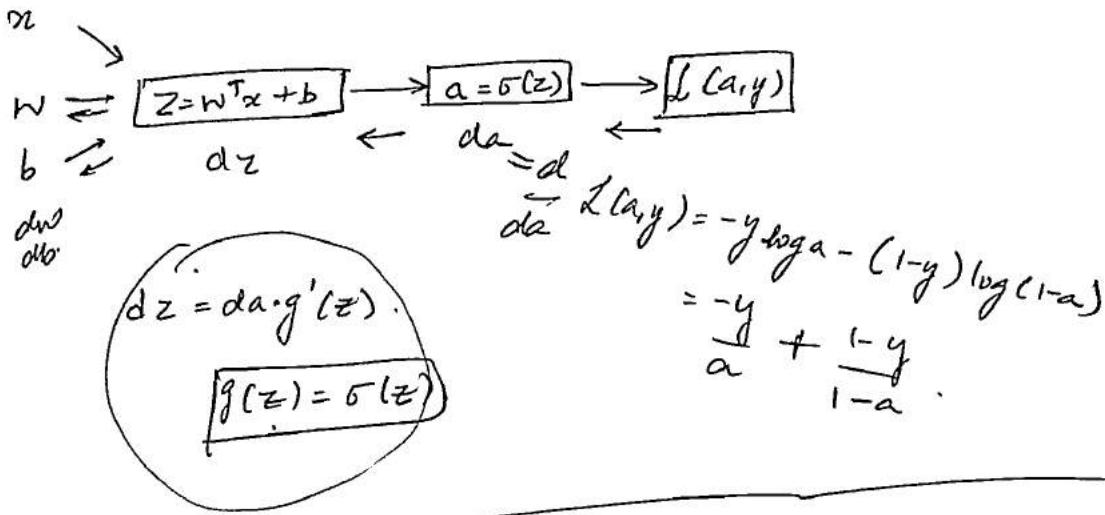
Backward Propagation

$$\begin{aligned} dz^{[2]} &= A^{[2]} - y \\ dw^{[2]} &= \frac{1}{m} dz^{[2]} A^{[1]T}, \\ db^{[2]} &= \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims=True}). \\ dz^{[1]} &= W^{[2]T} dz^{[2]} + g^{[1]'}(z^{[1]}) \end{aligned}$$

element wise product

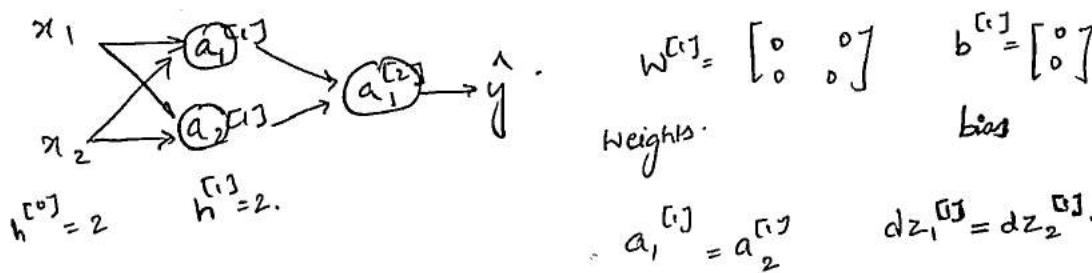
$$\begin{aligned} dW^{[i]} &= \frac{1}{m} dz^{[i]} X^T \\ db^{[i]} &= \frac{1}{m} \text{np.sum}(dz^{[i]}, \text{axis}=1, \text{keepdims=True}), \end{aligned}$$

Logistic regression



Random Initialization

* What happens if you initialize weights to zero?



dynamite:-

$$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \quad w^{[1]} = w^{[2]} - \alpha dw. \quad w^{[1]} = \begin{bmatrix} \dots \\ \dots \end{bmatrix}.$$

∴ if you initialize, with zero, w_1 after every iteration, will have the first row equal to the second row. ∵ They all are computing the same thing.

Initialize Randomly :-

$\therefore w^{(1)} = \text{np.random.rand}(2, 2) * 0.01$ Initialize to very small random values.

$$b^{[1]} = \text{np.zero}((z_{11}))$$

$$w^{(2)} = \dots$$

$$b^{(2)} = 0$$

W is random, so no longer there is a problem of symmetry breaking.

- ① Initialize parameters
 - ② predict w/ forward prop.
 - ③ Gradient Descent w/ Backward prop.

If it is very large, w will be large, Correspondingly $Z^{(1)} = W^{(1)} x + b^{(1)}$
 $a^{(1)} = g^{(1)}(Z^{(1)})$.

λ will be very large or small  the slope of the gradient will be very small so Gradient Descent will be very small & so leaving slow

LabPlanar Data - Classification with one hidden layer v6C.

np.random.seed(1).

2-class dataset

Mathematically,

For one example,

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

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

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

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

$$y_{\text{prediction}} = \begin{cases} 1 & \text{if } a^{[2](i)} > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

$$J = -\frac{1}{m} \sum_{i=0}^m \left(y^{(i)} \log(a^{[2](i)}) + (1-y^{(i)}) \log(1-a^{[2](i)}) \right)$$

General Net⁴ to build NN

① Define NN structure
(# of inp units, # of hidden units, etc.)

② Initialize the Model parameters

③ Loop:
- Implement forward propagation
- Compute loss
- Implement backward prop. to get gradients
- Update parameters.
(gradient descent)

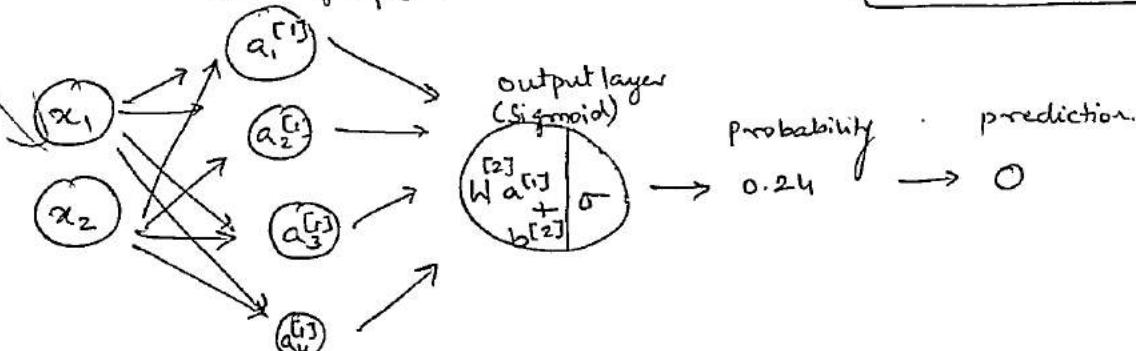
④ Define NN structure.

Define 3 variables:

↳ n_x : the size of the input layer↳ n_h : the size of the hidden layer (set this to 4)↳ n_y : the size of the output layer⑤ Use shapes of X & Y to find n_x and n_y

Also, hard code the hidden layer size to be 4.

hidden layer of size 4



Sigmoid()
np.tanh()

$$J = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(a^{[2]C(i)}) + (1-y^{(i)}) \log(1-a^{[2]C(i)}))$$

Implement `compute_cost()` to compute the value of the cost J .

$$-\sum_{i=0}^m y^{(i)} \log(a^{[2]C(i)}):$$

$$\text{log_probs} = \text{np.multiply}(\text{np.log}(A2), Y)$$

$$\text{cost} = -\text{np.sum}(\text{log_probs}).$$

Hardest Part in Deep Learning is Backward Propagation

17

Summary of Gradient Descent:-

~~$$dz^{[2]} = A^{[2]} - Y$$~~

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

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

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g'(Z^{[1]}).$$

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

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

①

o Functions Used

o o X, Y

* `layer_sizes(X, Y)`

* `initialize_parameters(n_x, n_h, n_y)`

* `forward_propagation(X, parameters)`

* `compute_cost(A2, Y, parameters)`

* `backward_propagation(parameters, cache, X, Y)`

* `update_parameters(parameters, grads, learning_rate = 1.2)`

Summary

Forward propagation

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

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

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

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

$$A^{[L]} = g^{[L]}(z^{[L]}) = \hat{y}$$

○

Implementation info

Cost function:

Compute Cost func to check if your model is actually learning. -

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

Backward Propagation:-

$$\text{d}W^{[L]} = \frac{\partial J}{\partial W^{[L]}} = \frac{1}{m} \text{d}z^{[L]} A^{[L-1]T}.$$

$$\text{d}b^{[L]} = \frac{\partial J}{\partial b^{[L]}} = \frac{1}{m} \sum_{i=1}^m \text{d}z^{[L]}(i).$$

$$\text{d}A^{[L-1]} = \frac{\partial L}{\partial A^{[L-1]}} = W^{[L]T} \text{d}z^{[L]}.$$

Backward Propagation

$$\text{d}z^{[L]} = A^{[L]} - y$$

$$\text{d}W^{[L]} = \frac{1}{m} \text{d}z^{[L]} A^{[L-1]T}$$

$$\text{d}b^{[L]} = \frac{1}{m} \text{np. sum}(\text{d}z^{[L]}, \text{axis}=1, \text{keep dims}=\text{True}).$$

$$\text{d}z^{[L-1]} = W^{[L]T} \cdot \text{d}z^{[L]} * g'^{[L]}(z^{[L]})$$

$$\text{d}z^{[1]} = \text{d}W^{[2]} \text{d}z^{[2]} * g'^{[1]}(z^{[1]}).$$

$$\text{d}W^{[1]} = \frac{1}{m} \text{d}z^{[1]} A^{[0]T}.$$

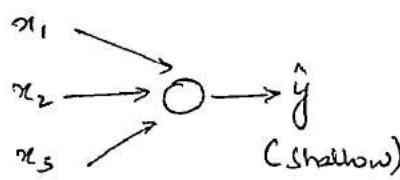
Note that $A^{[0]T}$ is another way to see previous page.

Wk4

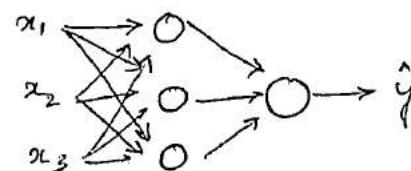
Deep Neural Networks - Deep L-layer Neural Network

* Deep Neural Network

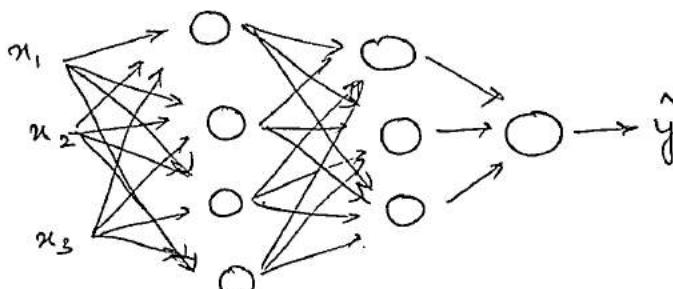
• 1 layer NN



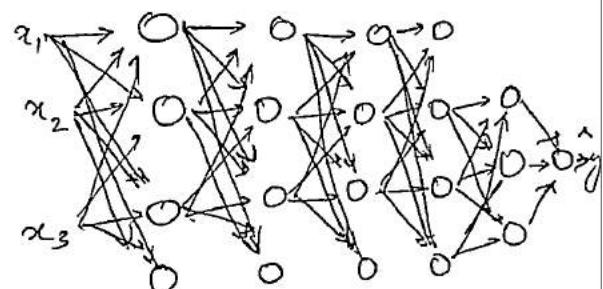
logistic regression "shallow"



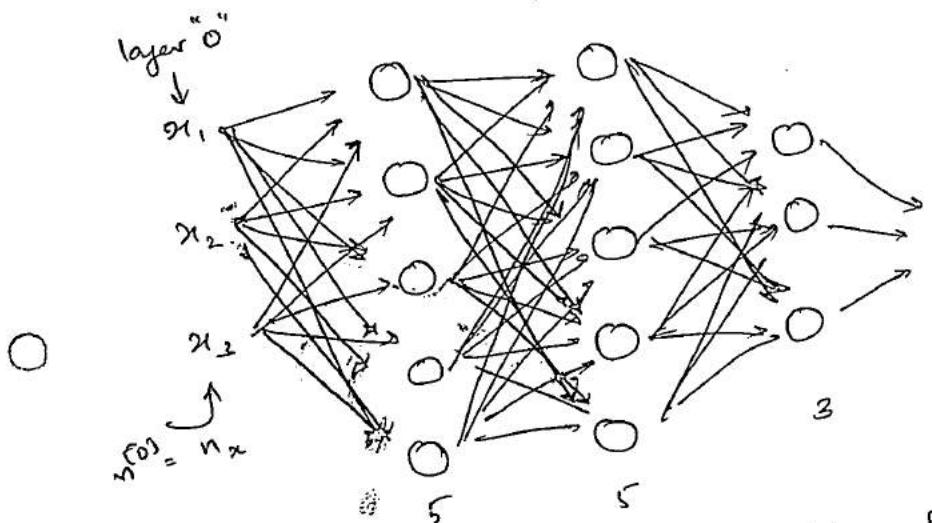
1 hidden layer.
 (2 Layer) NN



2 hidden layers.



"deep".
 5 - hidden layers.



4 layer NN

$$\text{O} \rightarrow \hat{y} = a^{[L]}$$

$L = 4$ (# layers); $n^{[l]} = \# \text{units in layer } l$. $n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3,$
 $n^{[4]} = 3, n^{[4]} = n^{[L]} = 1$

$$n^{[0]} = n_x = 3.$$

$a^{[l]}$ = activations in layer l .

$$a^{[l]} = g(z^{[l]}), \quad w^{[l]} = \text{weights for } z^{[l]}, \quad b^{[l]} = \text{biases } \dots$$

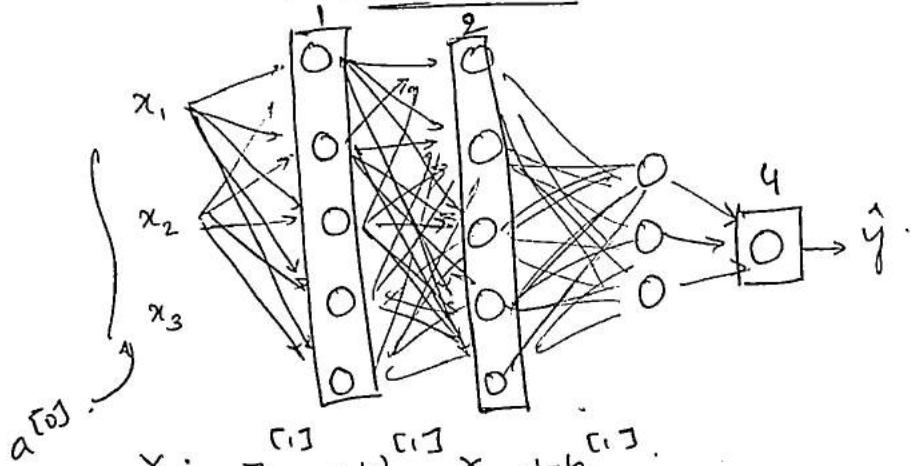
Activation index

finally,

$$X = a^{[0]}$$

input layer & also activation at index 0.

* Forward propagation in a deep network



$$x: z^{[1]} = \underbrace{W^{[1]} x}_{\text{effect the activations in layer } 1} + b^{[1]}.$$

activation fun^c

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

$$z^{[2]} = \underbrace{W^{[2]} a^{[1]}}_{\text{Parameters from [1] layer}} + b^{[2]}.$$

$$a^{[2]} = \underbrace{g^{[2]}(z^{[2]})}_{\text{activation fun^c$$

$$z^{[u]} = \underbrace{W^{[u]} a^{[u-1]}}_{\text{Parameters from [u-1] layer}} + \underbrace{b^{[u]}}_{\text{bias from [u] layer}}$$

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

Estimated Output \hat{y} is computed.

∴ Vectorized -

$$z^{[1]} = W^{[1]} \times \underbrace{A^{[0]}}_{\text{seems like a vector for } l=1 \dots 4} + b^{[1]}.$$

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

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

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

$$Y = g(z^{[4]}) = A^{[4]}.$$

Stacking columns
for m training examples.

$$\begin{bmatrix} z^{[2](1)} \\ z^{2} \\ \vdots \\ z^{[2](m)} \end{bmatrix} = \underline{\underline{z^{[2]}}}.$$

Act^c for layers
1, 2, 3, ...
Perf^c for Explicit for loops.

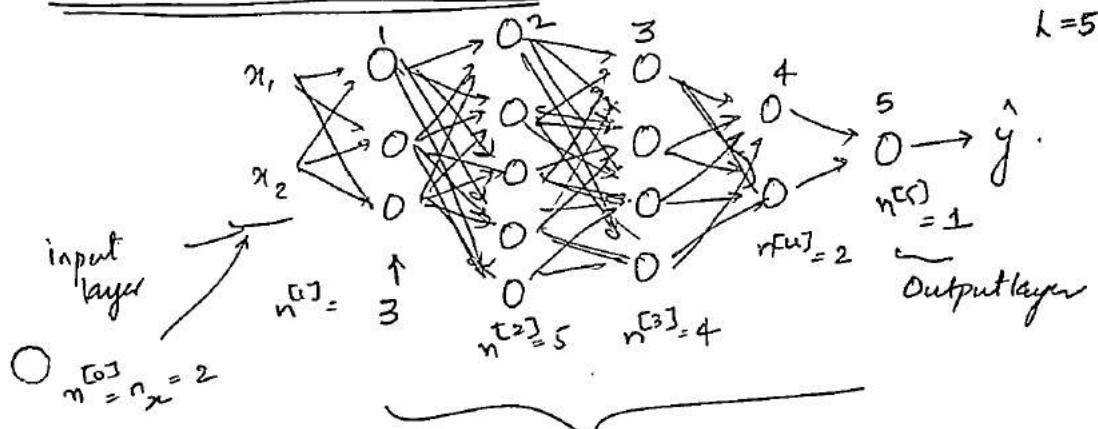
Correction for Video

$$a^{[e]} = g^{[e]}(z^{[e]}). \quad a, z \text{ have dimensions } (n^{[e]}, 1).$$

Wacht-3

Getting your matrix dimensions right

Parameters $\ln^{[e]}$ and $b^{[e]}$.



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

$$(3,4) \leftarrow (3,2)(2,1) + (3,1)$$

$$(n^{[r]},_1)(n^{[r]},_1 \underbrace{A(n^{[r]},_1)}_{n^{[r]},_1} + n^{[r]},_1).$$

$$\begin{bmatrix} : \\ : \\ : \end{bmatrix}_{3 \times 1} = \begin{pmatrix} : & : \\ : & : \\ : & : \end{pmatrix}_{3 \times 2} \begin{bmatrix} : \\ : \end{bmatrix}_{2 \times 1} +$$

From Matrix
Multiplication

Hidden layers

$$\omega^{[1]} : (n^{[1]}, n^{[0]}).$$

$$w^{[2]} : (5, 3) : (n^{[2]}, n^{[1]})$$

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

$$(5,1) \quad (5,3) \quad (3,1) + (5,1) \\ (n^2,1)$$

$$n^{[3]} \cdot s(4,5) \quad (n^{[3]}, n^{[2]})$$

$$w^{[4]}_{\text{opt}}(2,4), w^{[5]}_{\text{opt}}(1,2).$$

$$w^{[e]} = (n^{[e]}, n^{[e-1]}), \quad b^{[e]} = (n^{[e]}, 1)$$

Dimensions of your matrices.

for Back propagation:

$$dW^{[2]} : (n^{[1]}, n^{[1 \rightarrow 2]})$$

$$d b^{(e)} : (n^{(e)}, 1)$$

same as

$$\zeta^{[e]} = g^{[e]}(a^{[e]}).$$

Vectorized Implementation :-

Normally (Nonvector)

$$Z^{[l]} = W^{[l]} \cdot X + b^{[l]}$$

$(n^{[l]}, 1)$ $(n^{[l]}, n^{[0]})$ $(n^{[0]}, 1)$ $(n^{[l]}, 1)$

for Vectorized (Stacking them).

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

$$Z^{[l]} = W^{[l]} \cdot X + b^{[l]}$$

$(n^{[l]}, m)$ $(n^{[l]}, n^{[0]})$ $(n^{[0]}, m)$ $\underbrace{(n^{[l]}, 1)}$

through Python Broadcasting

$(n^{[l]}, m)$

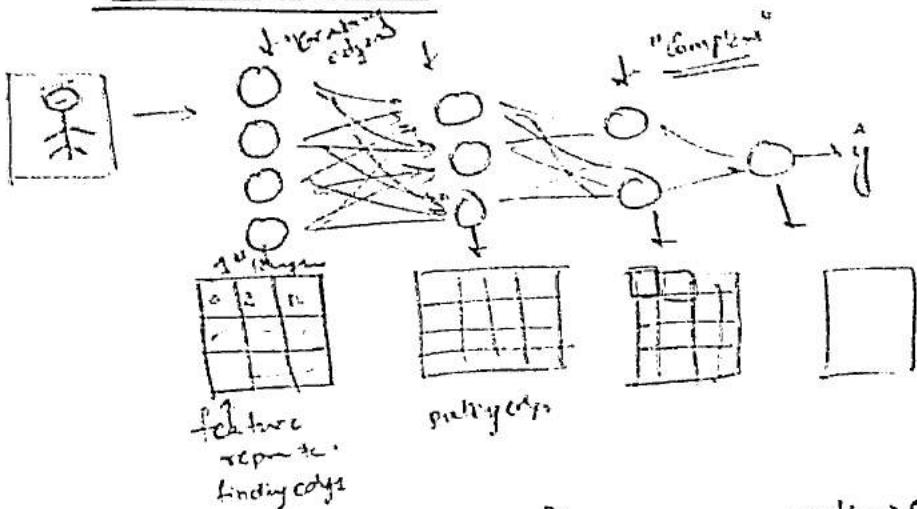
$$Z^{[l]}, a^{[l]} : (n^{[l]}, 1)$$

$$Z^{[l]}, A^{[l]} : (n^{[l]}, m)$$

capitalized (vectorized)
if $l=0 \rightarrow A^{[0]}=X = (n^{[0]}, m)$

$$dZ^{[l]}, dA^{[l]} : (n^{[l]}, m)$$

Intuition about deep representation



Audio \rightarrow low level \rightarrow Phonemes \rightarrow words \rightarrow sentence/phrases
 Audio waveform CAT

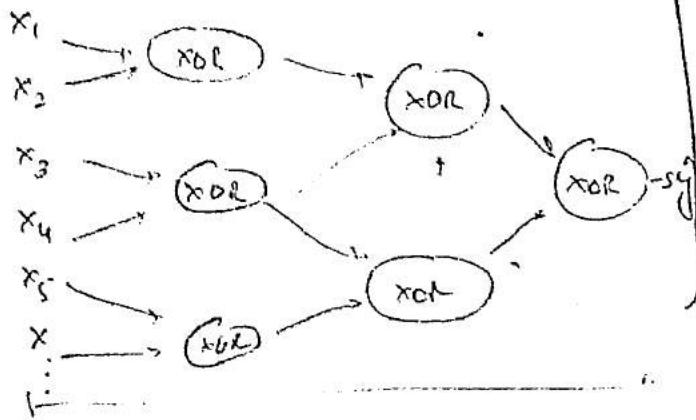
Circuit Theory and Deep Learning

Informally: There are functions you can compute with a "small" k-layer deep neural network that shallower networks require exponentially more hidden units to compute.

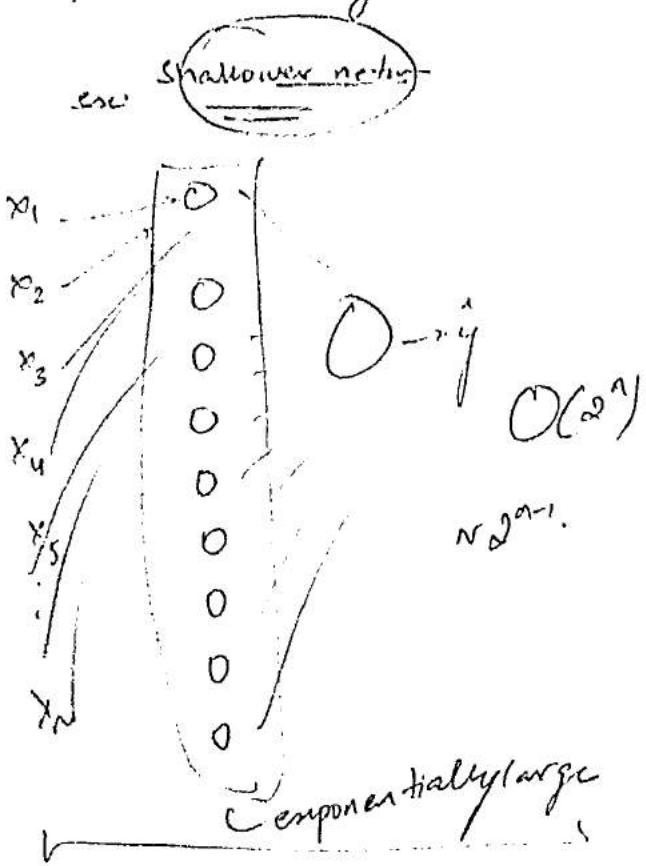
Deep Netw-

$$y = x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } \dots \text{ XOR } x_n.$$

$O(\log n)$



More efficient



Building blocks of Deep Neural Networks:

$$\begin{matrix} x_1 & 0 & 0 & 0 \\ x_2 & 0 & 0 & 0 \\ x_3 & 0 & 0 & 0 \\ x_4 & 0 & 0 & 0 \end{matrix}$$

$$0 \rightarrow \hat{y}$$

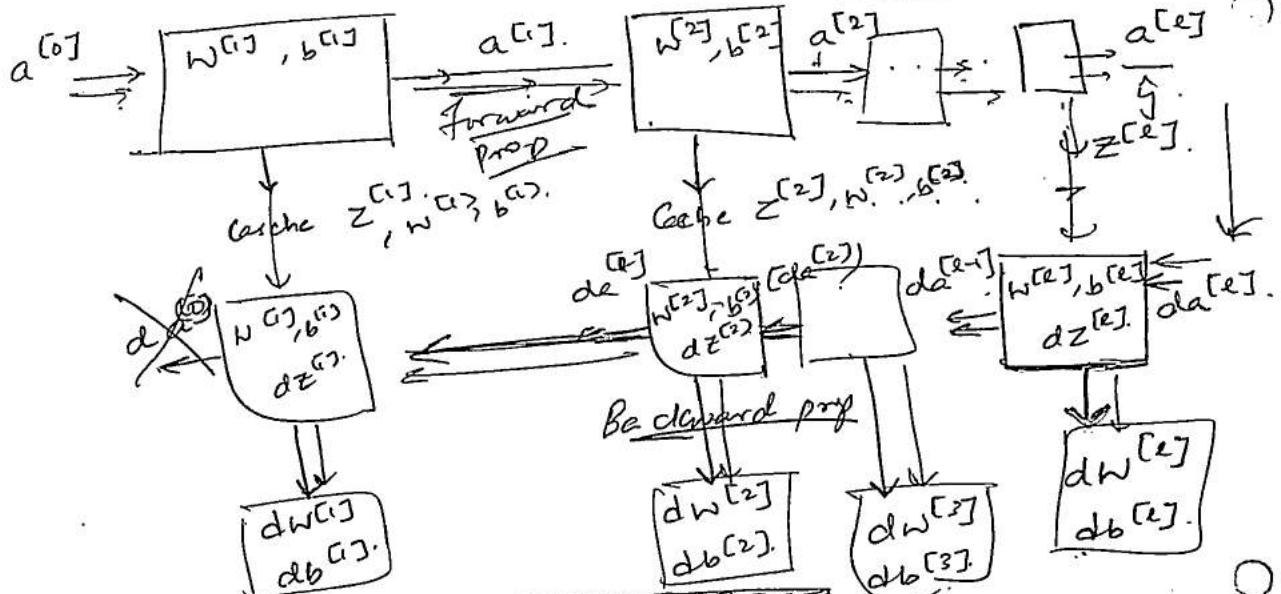
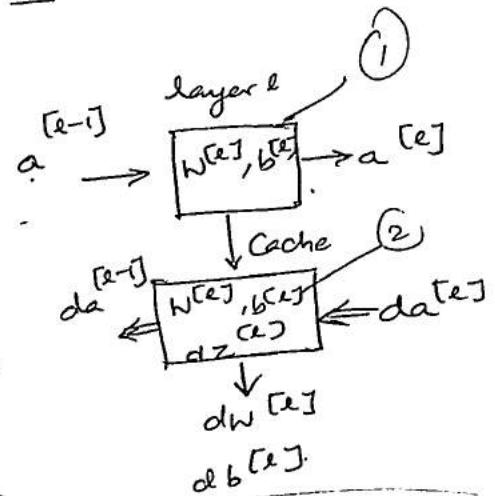
layer $l: w^{[e]}, b^{[e]}$

forward: input $a^{[e-1]}$, output $a^{[e]}$.

$$z^{[e]} = w^{[e]} a^{[e-1]} + b^{[e]}. \quad \text{Cache } z^{[e]}.$$

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

backward: input $d_a^{[e]}$, output $d_a^{[e-1]}$.
Cache $(z^{[e]})$



$$\boxed{\begin{aligned} w^{[e]} &:= w^{[e]} - \alpha d_w^{[e]} \\ b^{[e]} &:= b^{[e]} - \alpha d_b^{[e]}. \end{aligned}}$$

Correct Video

$$d_w^{[e]} = d_z^{[e]} * a^{[e-1]T}$$

Forward propagation for layer i

→ Input $a^{[e-i]}$.

→ Output $a^{[e]}$, Cache ($z^{[e]}$).

$$z^{[e]} = w^{[e]} \cdot a^{[e-i]} + b^{[e]}$$

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

Vectorized

$$z^{[e]} = w^{[e]} \cdot A^{[e-i]} + b^{[e]}$$

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

$$X = A^{[0]} \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \dots$$

Backward Propagation for layer i



→ Input $da^{[e]}$.

→ Output $da^{[e-i]}, dw^{[e]}, db^{[e]}$.

$$\frac{dz^{[e]}}{da^{[e]}} = g^{[e]}'(z^{[e]})$$

$$dw^{[e]} = dz^{[e]} * a^{[e-i]T}$$

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

$$da^{[e-i]} = w^{[e]T} \cdot dz^{[e]}$$

$$dz^{[e]} = w^{[e+1]T} \cdot dz^{[e+1]} * g^{[e]}'(z^{[e]})$$

Vectorized

$$dz^{[e]} = dA^{[e]} * g^{[e]}'(z^{[e]})$$

$$dw^{[e]} = \frac{1}{m} dz^{[e]} \cdot A^{[e-i]T}$$

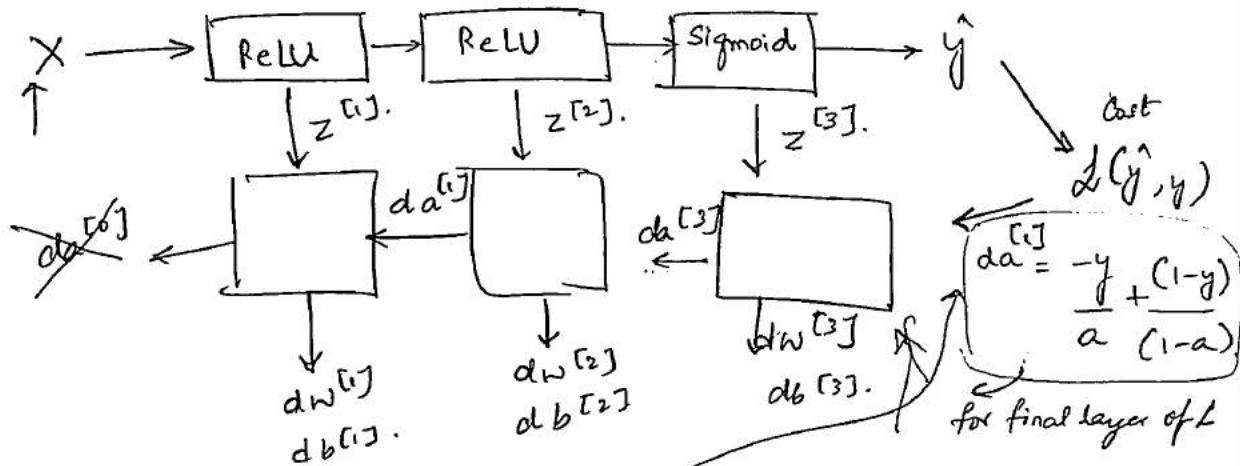
$$db^{[e]} = \frac{1}{m} \text{np.sum}(dz^{[e]}, \text{axis}=1)$$

$$dA^{[e-i]} = w^{[e]T} \cdot dz^{[e]}$$

(keepdims=True)

Backpropagation

SUMMARY



This is how you
Initialize the vectorized version of the
Back propagation -

$$dA^{[e]} = \left(\begin{array}{c} -\frac{\hat{y}^{(1)}}{a^{(1)}} + \frac{(1-\hat{y}^{(1)})}{(1-a^{(1)})} \\ \vdots \\ -\frac{\hat{y}^{(m)}}{a^{(m)}} + \frac{(1-\hat{y}^{(m)})}{(1-a^{(m)})} \end{array} \right)$$

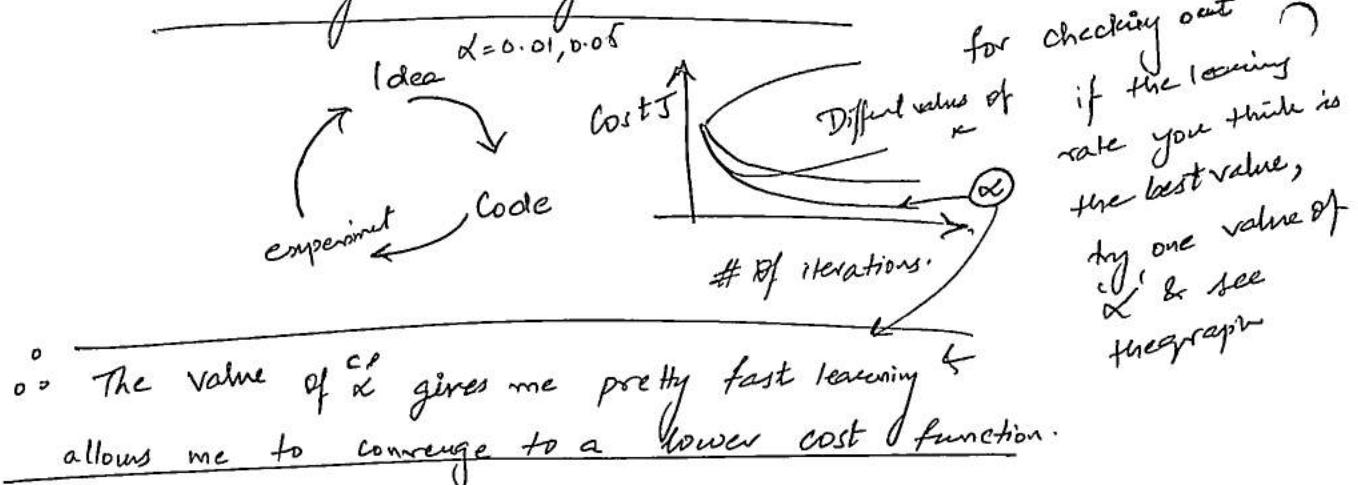
* Parameters vs Hyperparameters:

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots$

Hyperparameters:
learning rate α
iterations
hidden layers L
hidden units $n^{[1]}, n^{[2]}, \dots$
Choice of activation function.

Later: Momentum, mini batch size, regularizations...

Note * Applied deep learning is a very empirical process.



∴ The value of α gives me pretty fast learning & allows me to converge to a lower cost function.

// ① Try out many different values & check out which is best.
②

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

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

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

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

∴ \otimes element wise multiplication

$$dZ^{[1]} = W^{[2]} dZ^{[2]} * g'(z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} A^{[0]T}$$

Note that $A^{[0]T}$ is another way to denote the input features, which is also written as X^T . $db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims=True})$.