

Name:- Ashishay Jain
ashishay

PAGE NO.

DATE

This notebook contains the notes of:-

- DeepLearning.ai course 1 (Neural networks)
- " " " 2 (optimization, etc)
- " " " 3 (structuring ml projects)
- " " " 4 (CNN and computer vision, neural style, transfer)
+ (A section which further explains anchor boxes)
- Notes from the API project
(Syntax notes on pandas, numpy and matplotlib)
- Re-inforcement learning
- R programming
- RNN's from O'reilly book

0321448844
name: emp no:

Q1 Random forest algorithm.

- It is an ensemble method (based on divide and conquer approach) of decision trees generated on a randomly split dataset.
- This collection of decision tree classifiers is also known as forest.
- The individual decision trees are generated using an attribute selection indicators such as information gain, gain ratio, and gini index for each attribute.
- Each tree depends upon an independent random sample.
- In a classification problem, each tree votes and the most popular class is chosen as the final result.
- In case of regression, the average of all the tree outputs is considered as the final result.

The algorithm works in 4 steps:-

- 1) Select random samples from the given dataset.
- 2) Construct the decision tree for each sample and get a prediction result from each decision tree.

- 3.) Perform a vote for each predicted result.
- 4.) Select the prediction result with the most votes as final prediction.

~~Advantages:-~~

- Highly accurate as a lot of decision trees participate in the process.
- Does not suffer from overfitting as the average of all the predictions are taken which cancels out the biases.
- The algorithm can be used for both classification and regression problem.
- We can get relative feature importance.
- Random forests can also handle missing values. There are two ways to handle these:-
 - (i) Using median values to replace continuous variables.
 - (ii) Computing the proximity-weighted average of missing values.

~~Disadvantages:-~~

- Random forest is slow in generating predictions because it has multiple decision trees and it also has to perform

voting.

- The model is difficult to interpret as compared to decision tree, where you can easily make a decision by following the path in the tree.

Deep learning.

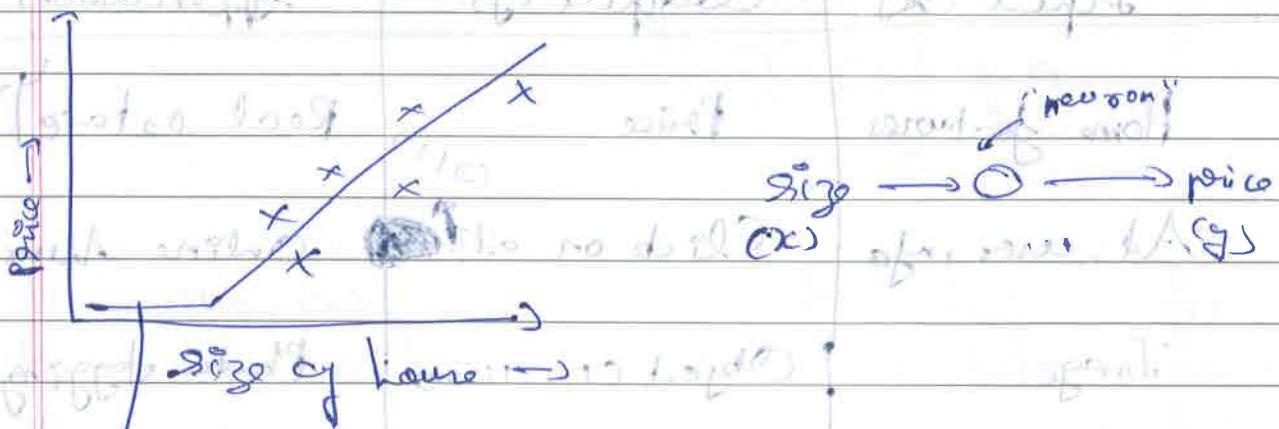
Part - 1

(1) Neural network and deep learning

(1.1) week (1)

(what is neural network)

→ Housing price prediction

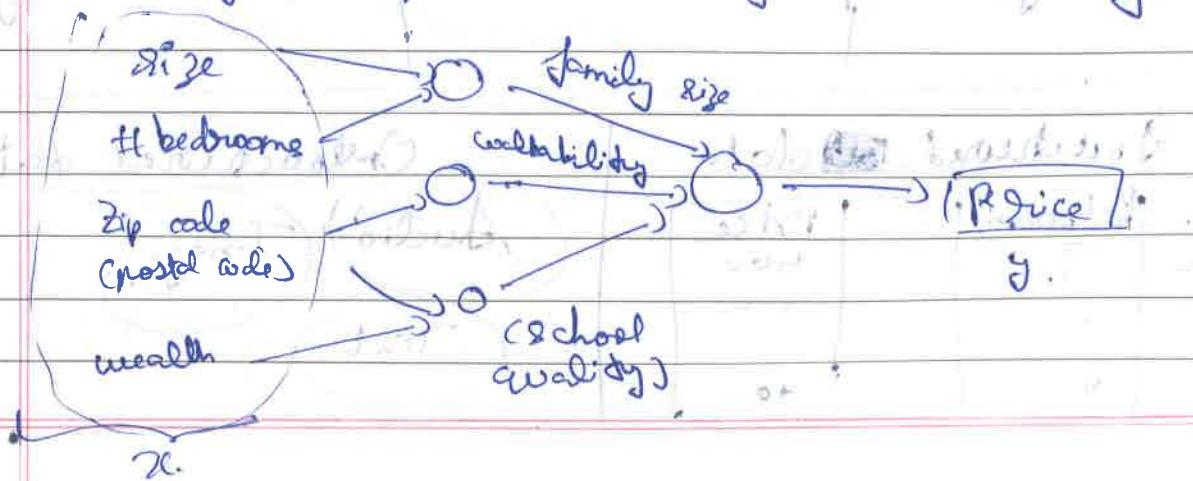


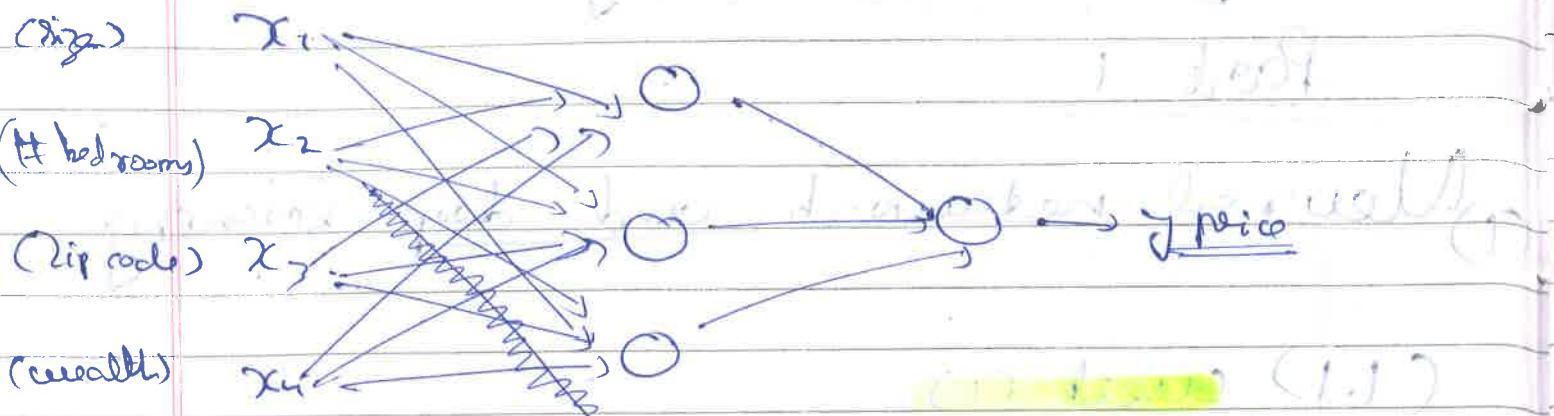
This ^{is} function is called Relu function

↳ Rectified linear unit

()

→ Housing price prediction using more than 1 features





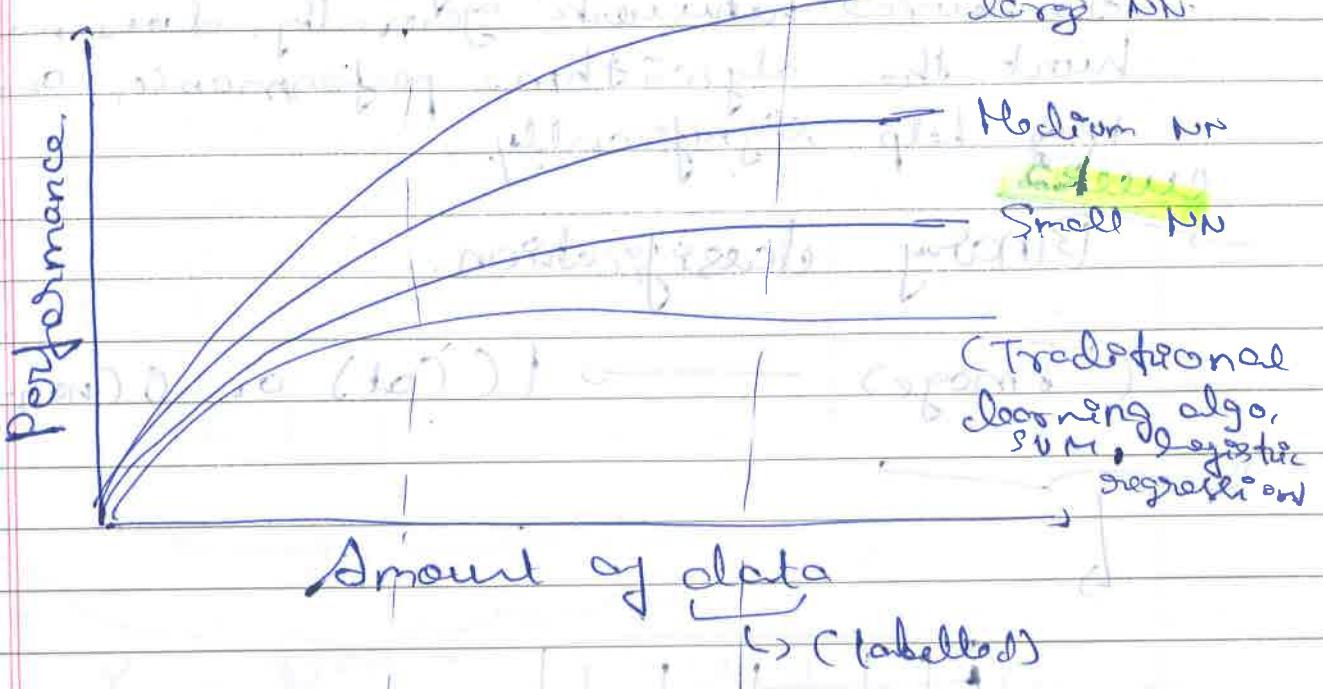
→ Supervised learning

Input (x_i)	Output (y)	Application
Home features	Price	Real estate
Ad. exec. info	Click on ad!	Online Advertising
Image	Object (1...1000)	Photo. tagging
Audio	Text document	Speech recognition
English	Chinese	Machine translation
Image, radar info	Position of other cars	Autonomous driving

Structured data	Unstructured data												
<table border="1"> <tr> <td>Size</td> <td>#bedrooms</td> <td>Price</td> </tr> <tr> <td>200</td> <td>3</td> <td>400</td> </tr> <tr> <td>1600</td> <td>3</td> <td>370</td> </tr> <tr> <td>3000</td> <td>4</td> <td>540</td> </tr> </table>	Size	#bedrooms	Price	200	3	400	1600	3	370	3000	4	540	Audio, Image, Text
Size	#bedrooms	Price											
200	3	400											
1600	3	370											
3000	4	540											

→ Scale drives deep learning progress

large NN



↓
network

Small training

set → large NN may
not dominate small NN or
traditional learning algo's.

↑
More larger NN

even better consistently

Thus the reason that deep learning is taking off :-

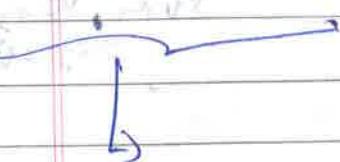
- (i) Deep learning has resulted in significant improvements in important applications such as online ~~chatbot~~ advertising, speech recognition and image recognition.
- (ii) We have access to a lot more computational power
- (iii) We have access to a lot more data.

And as a matter of fact increasing the neural network generally does not hurt the algorithms performance, and it may help significantly.

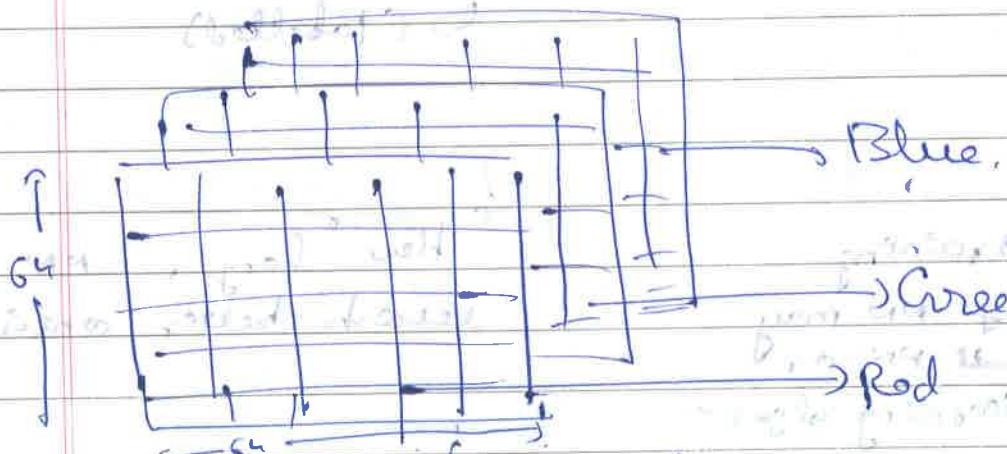
check

Binary classification.

C^2 (Image) \longrightarrow 1 (Cat) or 0 (Non-cat)



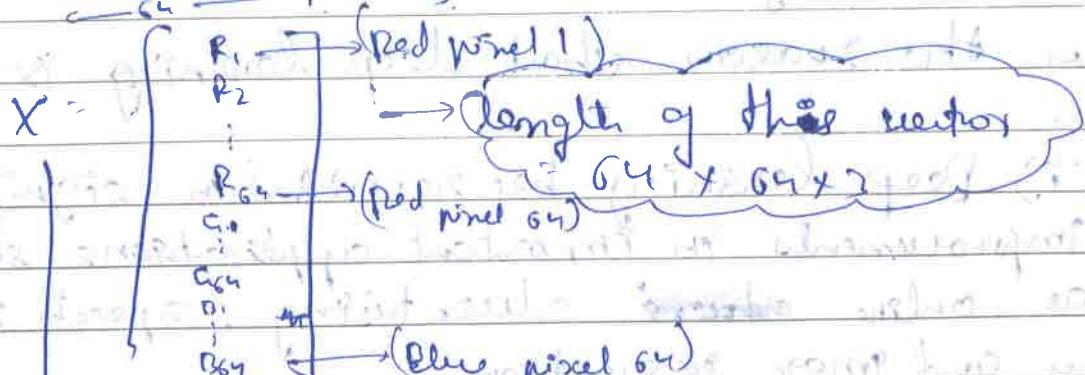
input to network.



Blue.

Green.

Red



\rightarrow This feature vector would be our input

It's used to (P.T.O)

→ Notation:-

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

in training example: $\{(x^{(0)}, y^{(0)}), (x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

$$X = \begin{bmatrix} | & | & | \\ x^{(0)} & x^{(1)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \quad n_x$$

$\leftarrow m \rightarrow$

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

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

$$y = [y^{(0)}, y^{(1)}, \dots, y^{(m)}]$$

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

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

→ logistic regression

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

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

parameters: $w \in \mathbb{R}^{n_x}$ and $b \in \mathbb{R}$

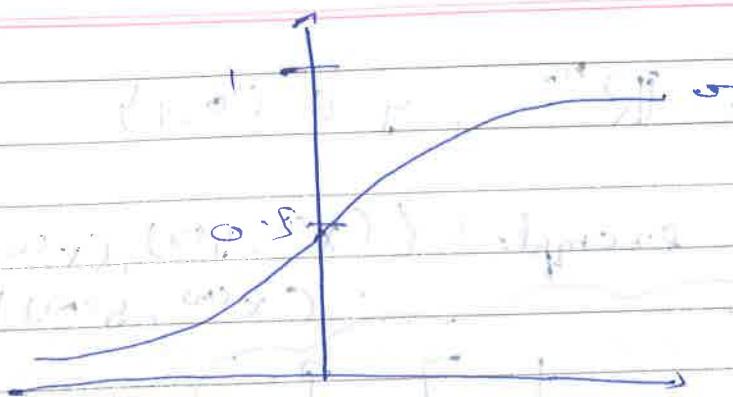
for linear regression we could have done:-

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

but for logistic regression, it must be

for logistic regression, it must be

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



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

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

If z is a large negative number

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

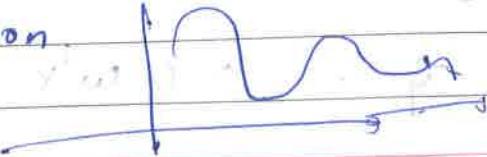
$$\rightarrow \hat{y} = \sigma(w^T x + b), \text{ where } \sigma(z) = \frac{1}{1 + e^{-z}}$$

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

loss (error) function:
we could have used

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

but we don't use that because it makes
gradient descent not work well as it is
a non-convex function.



thus we use the following loss function:-

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

If $y=1$: $L(\hat{y}, y) = -\log \hat{y}$ ← (want $\log \hat{y}$ large)
 If $y=0$: $L(\hat{y}, y) = -\log (1-\hat{y})$ ← (want $\log (1-\hat{y})$ large = want \hat{y} small)

→ Cost function:

$$\begin{aligned} J(w, b) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\ &= \frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right] \end{aligned}$$

→ The loss function computes the error for a single training example; the cost function is the average of the loss functions of the entire training set.

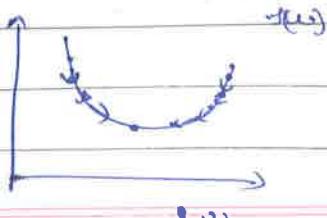
→ We want to find w, b to minimize $J(w, b)$

Gradient descent:-

Repeat {

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

}

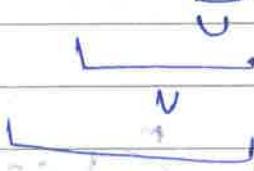


$$w := w - \alpha \frac{\partial J(c, w, b)}{\partial w}$$

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

→ Computation graph

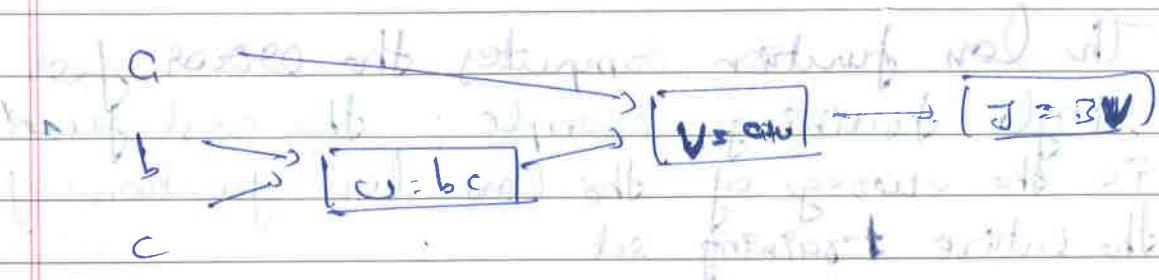
$$J(c, a, b, v) = 3(c + b)$$



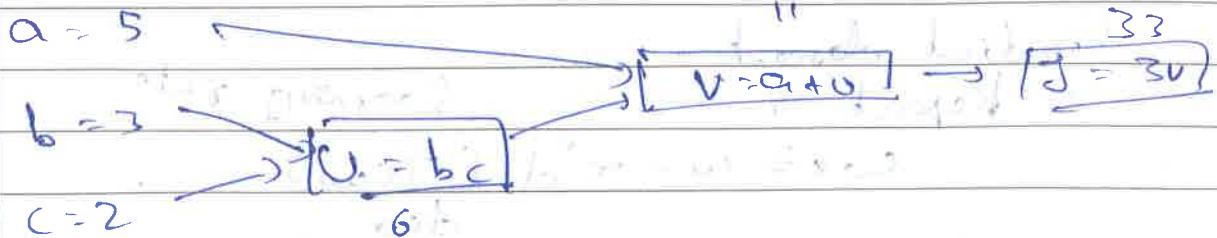
$$u = bc$$

$$N = a + b$$

$$J = 3u$$



→ Computing derivatives



$$J = 3N$$

$$V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

PAGE NO.

DATE

$$\text{thus } \frac{dJ}{dV} = 3$$

$$a = 5 \rightarrow 5.001$$

$$V = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003 \rightarrow \text{thus } \frac{dJ}{da} = 3$$

$$\frac{dJ}{da} = \left(\frac{dJ}{dV} \right) \cdot \left(\frac{dV}{da} \right)$$

The $\boxed{\Delta}$ amount by which J changes when a changes by Δa is the product by which J changes when we change V and Δ (by which V changes when we change a)

now from now on Δ will denote

Δ (Final output variable) as "divariable" (divariable)

$$a = 5$$

$$\Delta a = 3$$

$$b = 3$$

$$\Delta b = 6$$

$$c = 2$$

$$\Delta c = 9$$

Here we are calculating derivatives from right to left

→ logistic regression recap:

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

$$g = c = \sigma(z)$$

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

→ Logistic regression derivatives:

$$\begin{aligned}
 & \text{Inputs: } x_1, x_2, \dots \\
 & \text{Weights: } w_1, w_2, \dots, w_n, b \\
 & \text{Bias: } b \\
 & \text{Output: } z = w_1 x_1 + w_2 x_2 + \dots + b \\
 & \text{Sigmoid Function: } a = \sigma(z) = \frac{1}{1 + e^{-z}} \\
 & \text{Loss Function: } L(a, y) = -y \log a - (1-y) \log(1-a) \\
 & \frac{\partial L}{\partial w_1} = \frac{\partial L(a, y)}{\partial a} \cdot \frac{\partial a}{\partial w_1} = \frac{\partial L(a, y)}{\partial a} \cdot x_1 \\
 & \frac{\partial L}{\partial w_2} = \frac{\partial L(a, y)}{\partial a} \cdot \frac{\partial a}{\partial w_2} = \frac{\partial L(a, y)}{\partial a} \cdot x_2 \\
 & \frac{\partial L}{\partial b} = \frac{\partial L(a, y)}{\partial a} \cdot \frac{\partial a}{\partial b} = \frac{\partial L(a, y)}{\partial a} \cdot 1 \\
 & \frac{\partial L}{\partial a} = -y \cdot \frac{1}{a} + (1-y) \cdot \frac{1}{1-a} \\
 & \frac{\partial a}{\partial z} = \sigma'(z) = a(1-a) \\
 & \frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = -y \cdot a(1-a) + (1-y) \cdot a(1-a) \\
 & \frac{\partial L}{\partial z} = -y + y \cdot a(1-a) + (1-y) \cdot a(1-a) \\
 & \frac{\partial L}{\partial z} = -y + a - a \\
 & \frac{\partial L}{\partial z} = \underline{\underline{a - y}}
 \end{aligned}$$

$$\omega_1 := \omega_1 - \alpha \cdot d\omega_1 \quad \text{and update } \omega_1$$

$$\omega_2 := \omega_2 - \alpha \cdot d\omega_2 \quad \text{and update } \omega_2$$

$$b := b - \alpha \cdot db \quad \text{and update } b$$

→ Logistic regression on m examples.

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

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

so let's initialize

$$J = 0, \quad d\omega_1 = 0, \quad d\omega_2 = 0, \quad db = 0$$

For $i = 1$ to m

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

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

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

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

$$d\omega_1 = d\omega_1 + x_{1(i)} dz^{(i)}$$

$$d\omega_2 = d\omega_2 + x_{2(i)} dz^{(i)}$$

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

$$J = J/m$$

$$d\omega_1 = d\omega_1/m \quad \text{if } \alpha \text{ is large enough}$$

$$d\omega_2 = d\omega_2/m$$

$$db = db/m$$

$$\omega_1 := \omega_1 - \alpha \cdot d\omega_1$$

$$\omega_2 := \omega_2 - \alpha \cdot d\omega_2$$

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

After this we could be able to change ω_1, ω_2, b once. So we see that we have to use for loops & this would reduce our efficiency.

hence we ~~can~~ use vectorization to eliminate the use of for loops & increase the efficiency.

→ Vectorization:-

$$w = \begin{bmatrix} i \\ j \\ l \end{bmatrix} \quad x = \begin{bmatrix} i \\ j \\ l \end{bmatrix} \quad z = w^T x$$

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

non-vectorized:-

$$z = 0$$

for i in range(0, N):
 $z = z + w[i] * x[i]$

$$z = z + b$$

Vectorized

$$z = np.dot(w, x) + b$$

$$\curvearrowright w^T x$$

→ Python code for this in jupyter notebook:-

import numpy

a = np.array([1, 2, 3, 4])

print(a)

$\curvearrowright [1 2 3 4]$

import time

c = np.random.rand(1000000)

b = np.random.rand(1000000)

```

tic = time.time()
c = np.dot(a, b)
toc = time.time()
print ("Vectorized version: " + str(1000 * (toc - tic)) + " ms")

```

$c = 0$

```

tic = time.time()
for i in range (1000000):
    c = c + a[i] * b[i]
toc = time.time()
print ("For loop: " + str(1000 * (toc - tic)) + " ms")

```

Vectorized version: 1.505 ms

For loop : 481.710 ms

Thus we see that vectorized version is around 300 times faster.

The thing is that when we do vectorization and built-in libraries they use the parallel lines in GPU's and CPU's much more efficiently than for loop.

Whenever possible avoid explicit for loops.

→ Vectors and matrix valued functions.

Let $v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$

import numpy as np
 $v = np.exp(u) \rightarrow v = \begin{bmatrix} e^{u_1} \\ e^{u_2} \\ \vdots \\ e^{u_n} \end{bmatrix}$

$u = np.log(v) \rightarrow u = \begin{bmatrix} \log v_1 \\ \log v_2 \\ \vdots \\ \log v_n \end{bmatrix}$

$v = np.abs(u) \rightarrow v = \begin{bmatrix} \text{abs}(u_1) \\ \text{abs}(u_2) \\ \vdots \\ \text{abs}(u_n) \end{bmatrix}$

$u = np.array([1, 2, 3, 4, 5])$
 $v = np.maximum(u, 0) \rightarrow$ (between the maximum element
 elements of both the arrays)

$u \neq 2 \rightarrow v = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}$

→ Logistic regression derivatives.

$J = 0$ $dw = np.zeros(n, 1)$, $db = 0$

for $i = 1$ to m :

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

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

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

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

$$dw = dw + x^{(i)} dZ^{(i)}$$

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

$$J = J/m, dw = dw/m$$

$$db = db/m$$

→ Vectorized Logistic regression:-

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

$$\begin{aligned} [z^{(1)} \ z^{(2)} \dots z^{(m)}] &= \omega^T X + [b \dots b] \\ &= [\underbrace{\omega^T x^{(1)} + b}_{z^{(1)}} \ \underbrace{\omega^T x^{(2)} + b}_{z^{(2)}} \ \underbrace{\omega^T x^{(3)} + b}_{z^{(3)}} \dots \underbrace{\omega^T x^{(m)} + b}_{z^{(m)}}] \\ &= \underbrace{z}_{\text{Z}} \end{aligned}$$

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

python would automatically

convert b "real number" to (1xm) matrix & add each element b to np.dot(\omega^T, x). This is called "broadcasting".

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

now fill new we did

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

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

i

$$dZ = [dZ^{(1)} \ dZ^{(2)} \ dZ^{(3)} \dots dZ^{(m)}]$$

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

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

fill new we did:

{P.T.O}

till now we did :-

$$\rightarrow dw = 0$$

$$dw^1 = x^{(1)} dz^{(1)}$$

$$dw^2 = x^{(2)} dz^{(2)}$$

⋮

$$dw^m = x^{(m)} dz^{(m)}$$

$$dw^1 + \dots + dw^m = m$$

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

$$dw = \frac{1}{m} \left[\begin{matrix} x^{(1)} & \dots & x^{(m)} \end{matrix} \right] \left[\begin{matrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{matrix} \right]$$

$$= \frac{1}{m} [X^T dz^T + x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)}]$$

$$\rightarrow db = 0$$

$$db^1 = dz^{(1)}$$

$$db^2 = dz^{(2)}$$

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

$$db^1 + \dots + db^m = m$$

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

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

Then the code we will write is :-

for iter in range (0, 1000): -

$$\{ \begin{aligned} z &= w^T X + b \\ &= \text{np. dot}(w^T, x) + b \end{aligned}$$

$$A = \sigma(z)$$

$$dz = A - y$$

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

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

$$w := w - \alpha \cdot dw$$

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

→ Broadcasting Examples

Calories from carbs, proteins, Fat in 100 g of different foods:-

	Apples	Beef	Eggs	Potatoes
Carbs	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	175.0	99.0	0.9

We have to calculate % of calories from Carbs, protein, Fat.

By Import numpy as np.

$A = np.array([[56.0, 0.0, 4.4, 68.0], [1.2, 104.0, 52.0, 8.0], [1.8, 175.0, 99.0, 0.9]])$

$Cal = A.sum(axis=0)$

Now we will have $Cal = [59, 239, 155.4, 76.9]$

percentage = $100 * A / Cal$

Now percentage:-

$\frac{56}{59} * 100$	$\frac{0}{239} * 100$	$\frac{4.4}{155.4} * 100$	$\frac{68.0}{76.9} * 100$
$\frac{1.2}{59} * 100$	$\frac{104}{239} * 100$	$\frac{52.0}{155.4} * 100$	$\frac{8.0}{76.9} * 100$
$\frac{1.8}{59} * 100$	$\frac{175}{239} * 100$	$\frac{99.0}{155.4} * 100$	$\frac{0.9}{76.9} * 100$

This is the ans we require.

→ Some more based casting examples.

(i) $\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100$ \rightarrow $\begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$

→ $\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100$ \rightarrow $\begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$
 python would automatically convert to: $\begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix}$

(ii) $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}$

→ $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}$
 python would convert this to: $\begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$

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

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

→ $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}$
 python would convert this to

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

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

→ general principle:

$$(m, n) + (1, n) \rightarrow (m, n)$$

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

→

Tip:

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

$a.shape = (5,)$

"rank 0 array"

don't use

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

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

Try for some reason we are getting a rank 1 array then we should do:

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

→ logistic regression cost function:

so basically we have $\hat{y} = \sigma(w^T x + b)$:

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

We interpret \hat{y} as the probability that the image is a cat image given values x .
Thus

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

probability that $y=1$ given the values

the reason about \hat{y} is all we care of x .

$$\begin{cases} y=1 : P(y|X) = \hat{y} \\ y=0 : P(y|X) = 1-\hat{y} \end{cases}$$

- the probability of $y=1$ given the parameters X is \hat{y}
 → The probability of $y=0$ given the parameters X is $(1-\hat{y})$

Thus to summarise we could write:-

$$P(y|X) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

~~(i) If $y=1$ then the probability that $y=1$ given X would be \hat{y}~~

~~(ii) If $y=0$ then the probability that $y=0$ given X would be $(1-\hat{y})$~~

$$(i) \text{ If } y=1 \text{ then } P(y|X) = \hat{y}^1 (1-\hat{y})^{(1-1)} = \hat{y}$$

$$\text{hence } P(y=1|X) = \hat{y}$$

$$(ii) \text{ If } y=0 \text{ then } P(y|X) = \hat{y}^0 (1-\hat{y})^{(1-0)} = (1-\hat{y})$$

$$\text{hence } P(y=0|X) = (1-\hat{y})$$

thus we can write:-

$$P(y|X) = (\hat{y}^y) ((1-\hat{y})^{(1-y)})$$

now as \log is a strictly increasing function thus $\log(P(y|X))$ would

also indicate the probability.

$$\begin{aligned}
 \text{Thus } \log(p(y|x)) &= \log(g^{y|x} \cdot (1-g)^{1-y|x}) \\
 &= y \cdot \log g + (1-y) \cdot \log(1-g) \\
 &= -L(\hat{g}, y)
 \end{aligned}$$

thus as $\log(p(y|x)) \uparrow$, $L(\hat{g}, y) \downarrow$
 thus $p(y|x) \uparrow$, $L(\hat{g}, y) \downarrow$

and we need $p(y|x)$ to go up as we want more certainty and hence more probability thus we want $L(\hat{g}, y) \downarrow$, hence correctly termed as loss function.

→ Cost on m examples.

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

$$\begin{aligned}
 \Rightarrow \log p(\cdot) &= \sum_{i=1}^m \log p(y^{c_i} | x^{c_i}) \\
 &= -\sum_{i=1}^m L(\hat{g}^{c_i}, y^{c_i})
 \end{aligned}$$

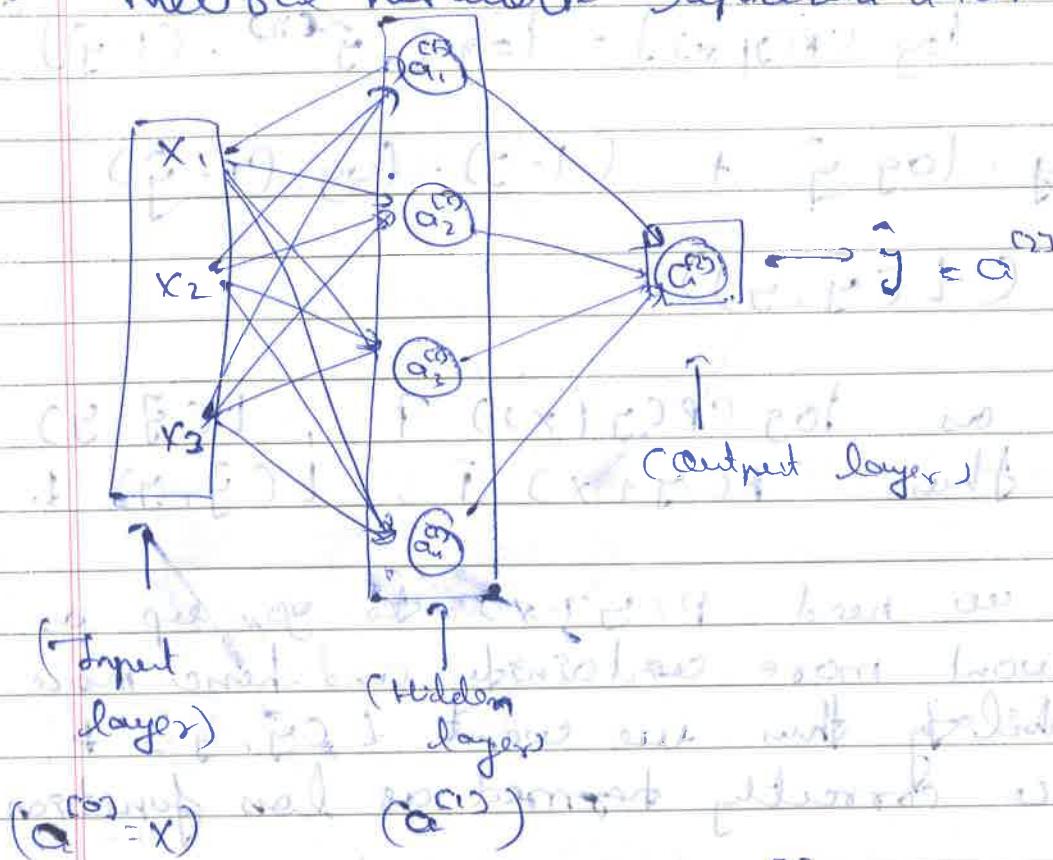
$$= -\sum_{i=1}^m L(\hat{g}^{c_i}, y^{c_i})$$

$$\text{Cost: } J(w, b) = \frac{1}{m} \left(\sum_{i=1}^m L(\hat{g}^{c_i}, y^{c_i}) \right)$$

(minimise)

week 3 (Shallow neural networks)

→ Neural network representation:

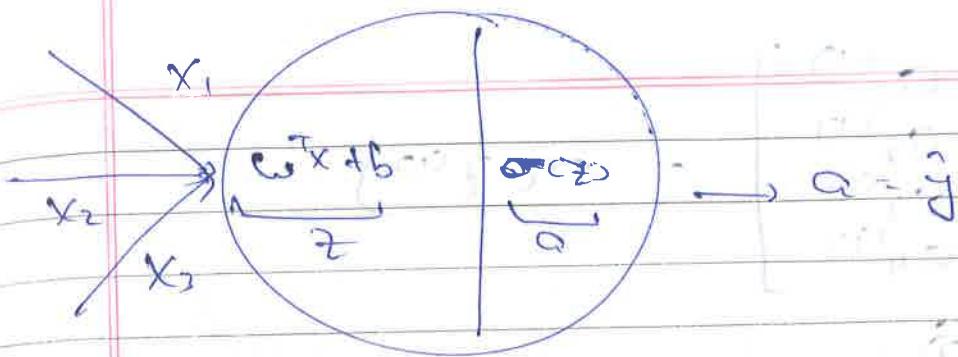


$$C^{(1)} = X \quad a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{bmatrix}$$

The above is a 2-layer neural network (The hidden layer & the output layer). we don't count the input layer as an official layer.

$$\begin{aligned} w^{(1)} &\rightarrow (4, 3) \rightarrow \text{dimensions} \\ b^{(1)} &\rightarrow (4, 1) \rightarrow \text{dimensions.} \end{aligned}$$

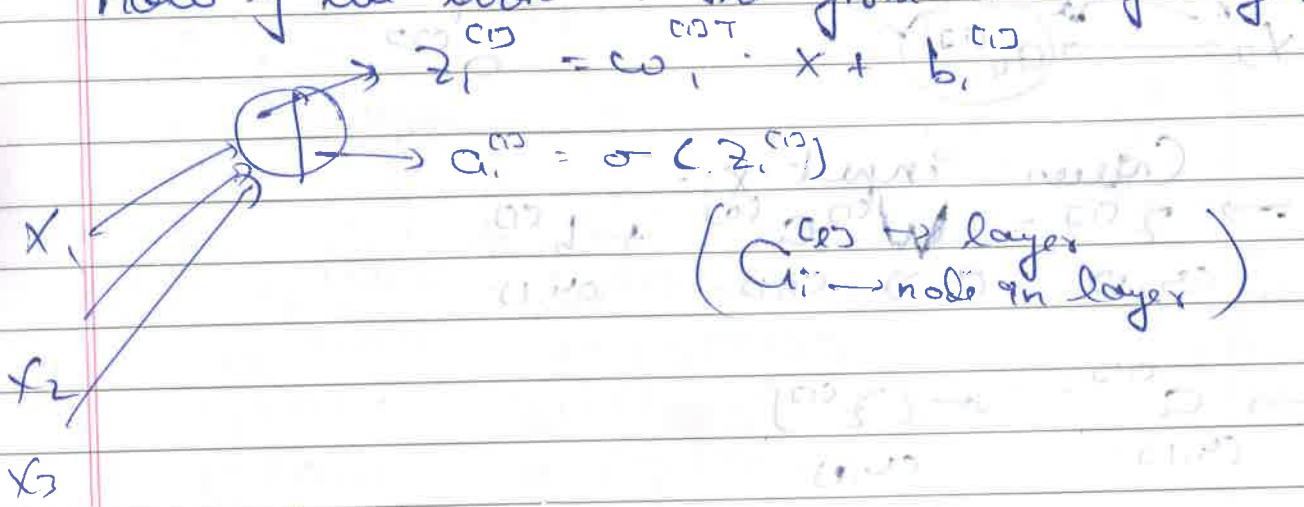
{P.T.O}



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

$$a = \sigma(z)$$

now if we look at the first node of layer 1:



Similarly:

$$z_1^{(1)} = w_1^{(1)T} x + b_1^{(1)}, \quad a_1^{(1)} = \sigma(z_1^{(1)})$$

$$z_2^{(1)} = w_2^{(1)T} x + b_2^{(1)}, \quad a_2^{(1)} = \sigma(z_2^{(1)})$$

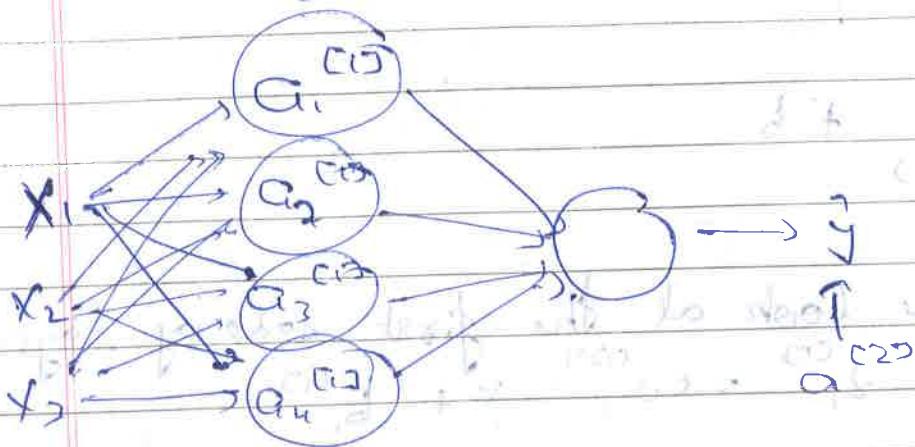
$$z_3^{(1)} = w_3^{(1)T} x + b_3^{(1)}, \quad a_3^{(1)} = \sigma(z_3^{(1)})$$

$$z_4^{(1)} = w_4^{(1)T} x + b_4^{(1)}, \quad a_4^{(1)} = \sigma(z_4^{(1)})$$

we could write this as

$$\begin{aligned}
 & \left[\begin{array}{c} w_1^{(1)T} \\ w_2^{(1)T} \\ w_3^{(1)T} \\ w_4^{(1)T} \end{array} \right] \left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right] + \left[\begin{array}{c} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{array} \right] = \left[\begin{array}{c} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{array} \right] \\
 & \left[\begin{array}{c} w_1^{(1)T} x + b_1^{(1)} \\ w_2^{(1)T} x + b_2^{(1)} \\ w_3^{(1)T} x + b_3^{(1)} \\ w_4^{(1)T} x + b_4^{(1)} \end{array} \right] = \left[\begin{array}{c} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{array} \right]
 \end{aligned}$$

$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{bmatrix} = \sigma(z^{(1)})$$



Given input x :

$$\rightarrow z^{(1)} = W^{(0)} a^{(0)} + b^{(0)}$$

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

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

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



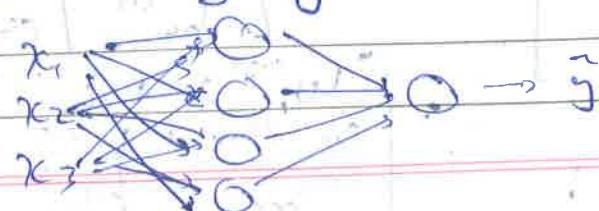
$$\rightarrow z^{(2)} = W^{(2)} a^{(1)} + b^{(2)}$$

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

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

$(1,1) \quad (1,1)$

\rightarrow Vectorizing across multiple examples.



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

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

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

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

$$X \longrightarrow a^{(2)} = \tilde{y}$$

$$X^{(1)} \longrightarrow a^{(2)(1)} = \tilde{y}^{(1)}$$

$$X^{(2)} \longrightarrow a^{(2)(2)} = \tilde{y}^{(2)}$$

$$X^{(m)} \longrightarrow a^{(2)(m)} = \tilde{y}^{(m)}$$

can we do example 1

into 1 layer 2 neurons of weight matrix

$$\begin{aligned} \rightarrow & \text{ for } i=1 \text{ to } m, \\ z^{(1)(i)} &= W^{(1)} X^{(i)} + b^{(1)} \\ a^{(1)(i)} &= \sigma(z^{(1)(i)}) \\ z^{(2)(i)} &= W^{(2)} a^{(1)(i)} + b^{(2)} \\ a^{(2)(i)} &= \sigma(z^{(2)(i)}) \end{aligned}$$

non
vectorized
implementation.

so we do:-

$$X = \begin{bmatrix} | & | & | & | \\ X^{(1)} & X^{(2)} & X^{(3)} & \dots & X^{(m)} \\ | & | & | & | & | \end{bmatrix} \quad W = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

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

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

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

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

$$Z^{(0)} = \begin{bmatrix} 1 & | & | & | & | & | \\ | & | & | & | & | & | \\ 1 & | & 0.001 & | & 0.002 & | & 0.003 \end{bmatrix} \xrightarrow{\text{Matrix Multiplication}} \begin{bmatrix} 1 & | & | & | & | & | \\ | & | & | & | & | & | \\ 1 & | & 0.001 & | & 0.002 & | & 0.003 \end{bmatrix}$$

$$A^{(0)} = \begin{bmatrix} 1 & | & | & | & | & | \\ | & | & | & | & | & | \\ a^{(0)01} & | & a^{(0)02} & | & \dots & | & a^{(0)m} \end{bmatrix} \xrightarrow{\text{Matrix Multiplication}} \begin{bmatrix} 1 & | & | & | & | & | \\ | & | & | & | & | & | \\ a^{(0)01} & | & a^{(0)02} & | & \dots & | & a^{(0)m} \end{bmatrix}$$

→ Justification for vectorized implementation.

$$Z^{(0)(0)} = W^{(0)} X^{(0)} + b^{(0)}, Z^{(0)(0)} = W^{(0)} X^{(0)} + b^{(0)}$$

$$W^{(0)} = \begin{bmatrix} \dots & | & \dots & | & \dots & | & \dots \end{bmatrix} \quad X^{(0)} = \begin{bmatrix} \dots & | & \dots & | & \dots & | & \dots \end{bmatrix}$$

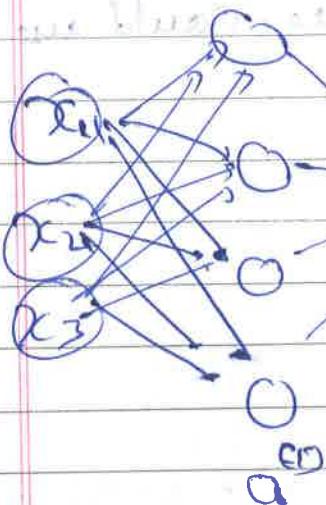
$$W^{(0)} X^{(0)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$W^{(0)} X^{(0)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$W^{(0)} X^{(0)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$W^{(0)} \begin{bmatrix} X^{(0)} & X^{(1)} & X^{(2)} & \dots \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

→ Recap of ~~conventional~~ vectorization errors
multiple examples:-



$$W^{(1)} = \begin{bmatrix} w_1^{(1)} & w_2^{(1)} & w_3^{(1)} & w_4^{(1)} \end{bmatrix}$$

$$x = \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \\ x_4^{(1)} \end{bmatrix}$$

$$x = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(m)} \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(m)} \\ x_3^{(1)} & x_3^{(2)} & \dots & x_3^{(m)} \\ x_4^{(1)} & x_4^{(2)} & \dots & x_4^{(m)} \end{bmatrix}$$

3x m

→ for $i \in (0, m)$:

$$Z_{(4, i)}^{(1)} = W_{(4, 3)}^{(1)} X_{(3, i)}^{(1)} + b_{(4, i)}^{(1)}$$

$$G_{(4, i)}^{(1)} = \sigma(Z_{(4, i)}^{(1)})^{(4, 1)}$$

$$Z_{(1, i)}^{(2)} = W_{(1, 4)}^{(2)} G_{(4, i)}^{(1)} + b_{(1, i)}^{(2)}$$

$$a_{(1, i)}^{(2)} = \sigma(Z_{(1, i)}^{(2)})_{(1, 1)}$$

vectorizing
error
i

$$Z_{(4, m)}^{(1)} = W_{(4, 3)}^{(1)} X_{(3, m)}^{(1)} + b_{(4, m)}^{(1)}$$

$$A_{(4, m)}^{(1)} = \sigma(Z_{(4, m)}^{(1)})$$

broadcasting

$$Z_{(1, m)}^{(2)} = W_{(1, 4)}^{(2)} \cdot A_{(4, m)}^{(1)} + b_{(1, m)}^{(2)}$$

broadcasting

$$A_{(1, m)}^{(2)} = \sigma(Z_{(1, m)}^{(2)})$$

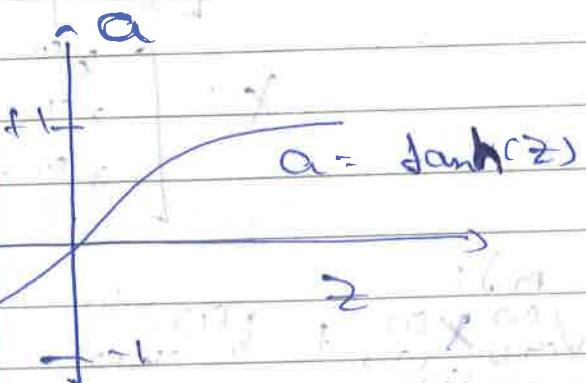
Note: $\tan^{-1}(x)$ and ~~$\tanh(x)$~~ $\tanh(x)$
are two different expressions.

PAGE No. /

→ Activation functions.

instead of using $\sigma(z)$ we could use $\tanh(z)$

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



thus if

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

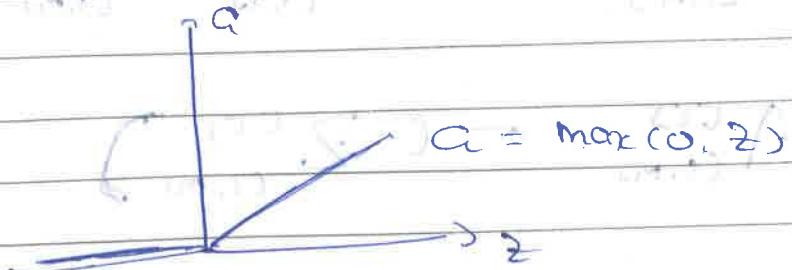
$$a^{(1)} = g^{(1)}(z^{(1)})$$

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

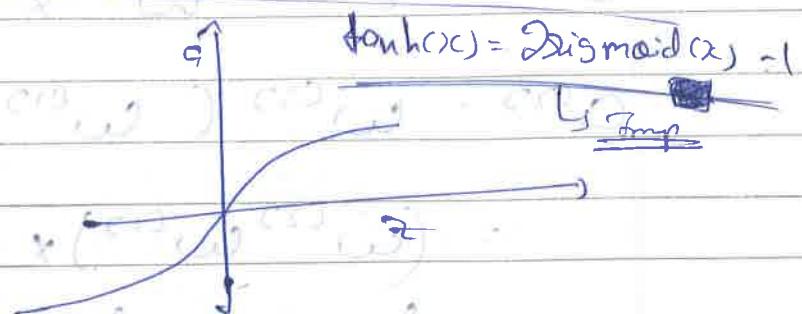
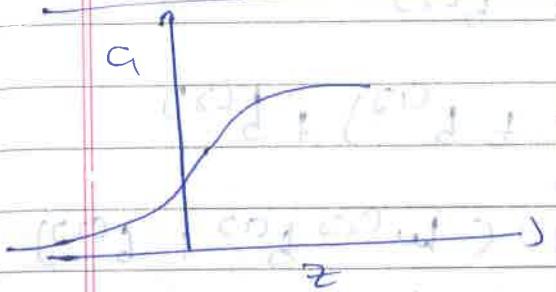
$$a^{(2)} = g^{(2)}(z^{(2)})$$

where g could be $\sigma()$ or $\tanh()$

→ we could also use ~~ReLU~~ ReLU function



→ Pros and cons of activation function



Sigmoid:

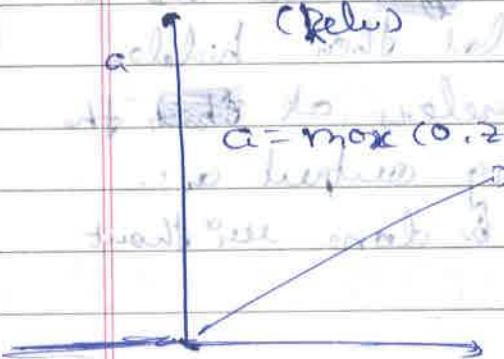
$$a = 1$$

(Slow and very less, only for binary classification)

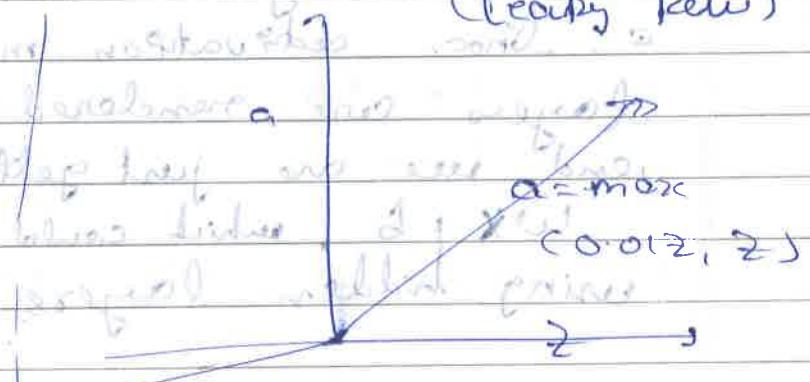
$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

(mostly all the time works better than sigmoid)

→



(by default uses this, mostly this works the best)



(Leaky ReLU)

→ Why does neural network need non-linear activation function.

$$\rightarrow \text{If we do } g(z) = z$$

"linear activation function"

$$z^{(0)} = w^{(0)}x + b^{(0)}$$

$$a^{(0)} = z^{(0)}$$

$$z^{(1)} = w^{(1)}a^{(0)} + b^{(1)}$$

$$a^{(1)} = z^{(1)}$$

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

$$a^{(2)} = z^{(2)} = w^{(2)}x + b^{(2)}$$

$$a^{(2)} = w^{(2)} (w^{(1)}x + b^{(1)}) + b^{(2)}$$

$$= (\underbrace{w^{(2)}w^{(1)}}_{w}, \underbrace{x}_{x}) + (\underbrace{w^{(2)}b^{(1)} + b^{(2)}}_{b'})$$

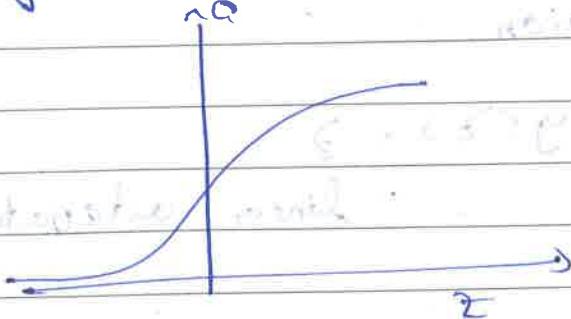
$$= w'x + b'$$

hence we get that if ~~we are~~ we are a linear activation model then hidden layers are rendered useless at ~~the~~ the end we are just getting output as $wx + b$, which could be done without using hidden layers.

We could use linear activation function sometimes too compute the final output at the last hidden layer ~~though~~ though.

→ Derivatives of activation function.

→ Sigmoid activation function



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

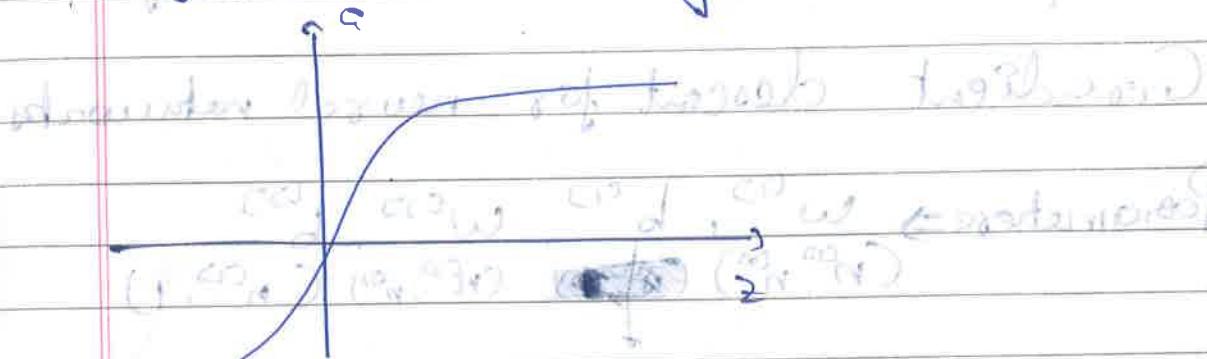
$$\frac{d(g(z))}{dz} = \text{slope of } g(z) \text{ at } z = g'(z)$$

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

$$= g(z) \cdot (1-g(z))$$

Thus if $g(z) = a$ then $g'(z) = a(1-a)$

\Rightarrow tan⁻¹ activation function.



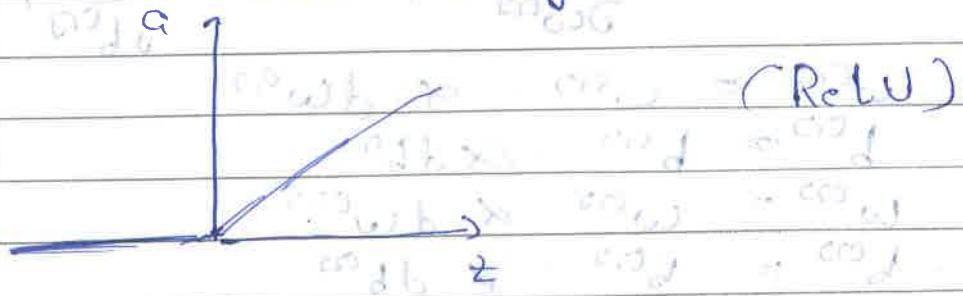
$$g(z) = \tan^{-1}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

$$= 1 - (g(z))^2$$

Then if $a = g(z)$ then $g'(z) = 1 - (g(z))^2 \equiv 1 - a^2$

\Rightarrow ReLU and leaky ReLU



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

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

Leaky Relu

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

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

→ Gradient descent for neural networks

Parameters $\rightarrow w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}$
 $(n^{(0)}, n^{(1)})$ $\xrightarrow{\text{---}}$ $(n^{(1)}, n^{(2)})$ $\xrightarrow{\text{---}}$ $(n^{(2)}, 1)$

cost function: $J(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)})$

$$= \frac{1}{m} \sum_{i=1}^m L(g_i, y_i)$$

Gradient descent:-

→ Repeat for

Compute predicts $(\hat{y}^{(i)}, i=1 \dots m)$
 $dJ = \frac{\partial J}{\partial w^{(1)}}, \frac{\partial J}{\partial b^{(1)}}, \dots, \frac{\partial J}{\partial b^{(2)}}$

$$w^{(1)} = w^{(1)} - \alpha \cdot dJ$$

$$b^{(1)} = b^{(1)} - \alpha \cdot dJ$$

$$w^{(2)} = w^{(2)} - \alpha \cdot dJ$$

$$b^{(2)} = b^{(2)} - \alpha \cdot dJ$$

→ formulas for computing derivatives
 Forward propagation:-

$$Z^{(0)} = W^{(0)} X + b^{(0)}$$

$$A^{(0)} = g^{(0)}(Z^{(0)})$$

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

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

back propagation:-

$$dZ^{(0)} = A^{(0)} - Y$$

$$dW^{(1)} = \frac{1}{m} (dZ^{(0)} A^{(0)T})$$

$$db^{(0)} = \frac{1}{m} \text{np.sum}(dZ^{(0)}, \text{axis}=1, \text{keepdims=True})$$

$$dZ^{(1)} = W^{(1)T} dZ^{(0)} \quad \underbrace{g^{(0)T}(Z^{(0)})}_{\text{element wise product}} \rightarrow (n^{(1)}, m)$$

$$dW^{(0)} = \frac{1}{m} dZ^{(1)} X^T$$

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

{P.T.O}

→ Backpropagation, intuition (optional)

→ Computing gradients:

logistic regression:

$$z = w^T x + b \quad a = g(z) = \frac{1}{1+e^{-z}}$$

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

$$dw = d_z \cdot x$$

$$db = d_z \quad \text{know that } a = g(z)$$

$$g(z) = \sigma(z) = a$$

$$\Rightarrow d(g(z)) = g'(z)$$

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

$$\Rightarrow g'(z) = \frac{da}{dz} \quad \text{means } \frac{da}{dz}$$

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

→ Neural network gradients:

$$z^{(0)} = w^{(0)} x + b^{(0)} \quad a^{(0)} = \sigma(z^{(0)})$$

$$d_z^{(0)} = w^{(0)} d_z^{(1)} + g^{(0)}(z^{(0)})$$

$$d_w^{(0)} = d_z^{(0)} \cdot x$$

$$d_b^{(0)} = d_z^{(0)}$$

$$z^{(1)} = w^{(1)} a^{(0)} + b^{(1)} \quad a^{(1)} = \sigma(z^{(1)})$$

$$d_z^{(1)} = w^{(1)} d_z^{(0)} + g^{(1)}(z^{(1)})$$

$$d_w^{(1)} = d_z^{(1)} \cdot a^{(0)}$$

$$d_b^{(1)} = d_z^{(1)}$$

$$z^{(2)} = w^{(2)} a^{(1)} + b^{(2)} \quad a^{(2)} = \sigma(z^{(2)})$$

$$d_z^{(2)} = w^{(2)} d_z^{(1)} + g^{(2)}(z^{(2)})$$

$$d_w^{(2)} = d_z^{(2)} \cdot a^{(1)}$$

$$d_b^{(2)} = d_z^{(2)}$$

$$z^{(3)} = w^{(3)} a^{(2)} + b^{(3)} \quad a^{(3)} = \sigma(z^{(3)})$$

$$d_z^{(3)} = w^{(3)} d_z^{(2)} + g^{(3)}(z^{(3)})$$

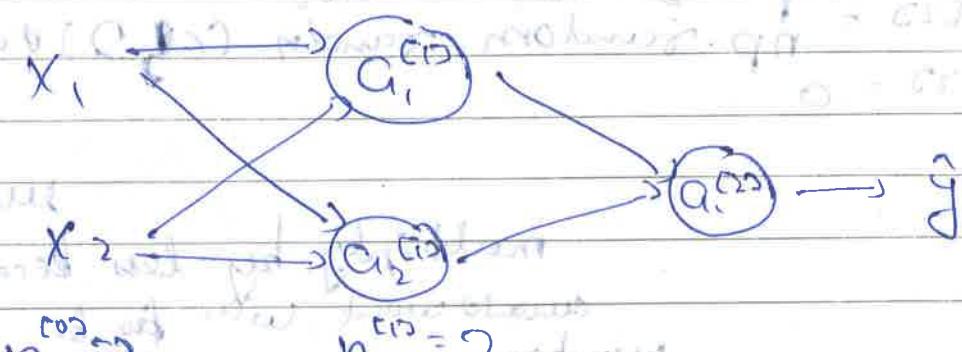
$$d_w^{(3)} = d_z^{(3)} \cdot a^{(2)}$$

$$d_b^{(3)} = d_z^{(3)}$$

$$[(a^{(2)}, g)]$$

→ Random Initialization.

What happens if you initialize weight to zero?



$$y, w^{(1)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Then we get $a_1^{(1)} = a_2^{(1)}$

$$w^{(2)} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$d2^{(1)} = d2^{(2)}$$

$$dw = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$w^{(2)} = w^{(2)} - \alpha \cdot dw$$

So basically we will always have dw in the form of $\begin{bmatrix} \alpha & \beta \\ \alpha & \beta \end{bmatrix}$

Always $a_1^{(1)} = a_2^{(2)}$ which would fail the purpose of hidden layers

so instead we should use:

$$w^{(1)} = np.random.randn(2, 2) * 0.01$$

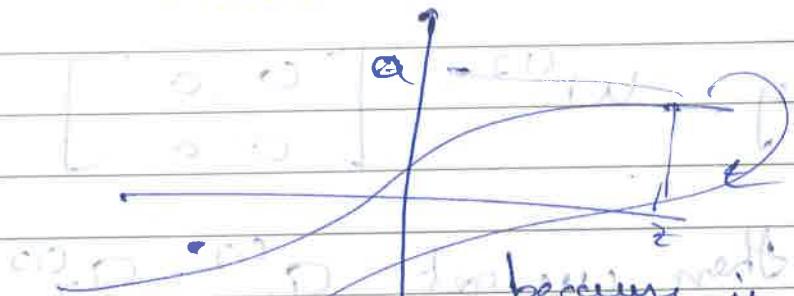
~~$b^{(1)} = np.zeros(2)$~~

$$b^{(1)} = np.zeros(1)$$

$$w^{(2)} = np.random.randn(3, 2) * 0.01$$

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

we would multiply by less because we would want $w^{(1)}$ to be small numbers

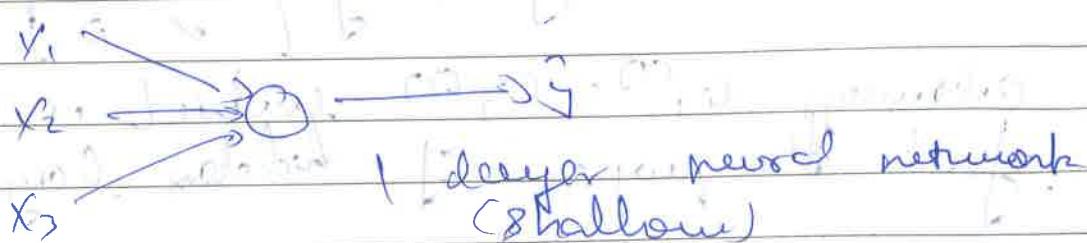


because if we small multiply by large then z would be large as $z = w^T x + b$
if z would be large then a would be around here

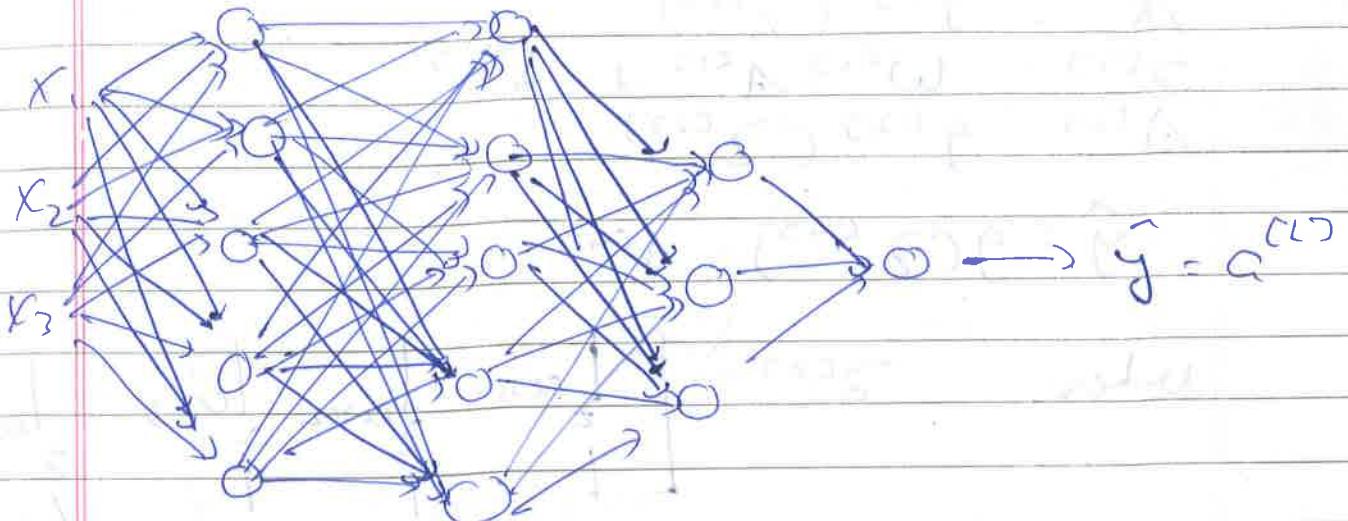
$$\text{then } da/dz \approx 0$$

which would make our algorithm very slow.

Week 5 (Deep L-layer neural network)



Deep neural network



$$l = 4 \quad (\# \text{ layers})$$

$a^{(l)}$ = # units in layer l .

$a^{(l)}$ = activations in layer l

$$a^{(l)} = g^{(l)}(z^{(l)}) \quad w^{(l)} = \text{weights for } z^{(l)}$$

$$n^{(0)} = 5, n^{(1)} = 5, n^{(2)} = 3, n^{(3)} = n^{(4)} = 1$$

$$n^{(0)} = n_x = 3$$

$$a^{(0)} = x$$

~~$$z^{(1)} = w^{(0)} x + b^{(0)}$$~~

$$z^{(1)} = w^{(0)} x + b^{(0)}$$

$$a^{(1)} = g^{(1)}(z^{(1)})$$

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

$$a^{(2)} = g^{(2)}(z^{(2)})$$

$$z^{(3)} = w^{(2)} a^{(2)} + b^{(2)}$$

$$a^{(3)} = g^{(3)}(z^{(3)})$$

Then do generalize:-

$$z^{(l)} = w^{(l-1)} a^{(l-1)} + b^{(l)}$$

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

Vectorized :-

$$z^{(1)} = w^{(1)} A^{(0)} + b^{(1)}$$

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

$$z^{(2)} = w^{(2)} A^{(1)} + b^{(2)}$$

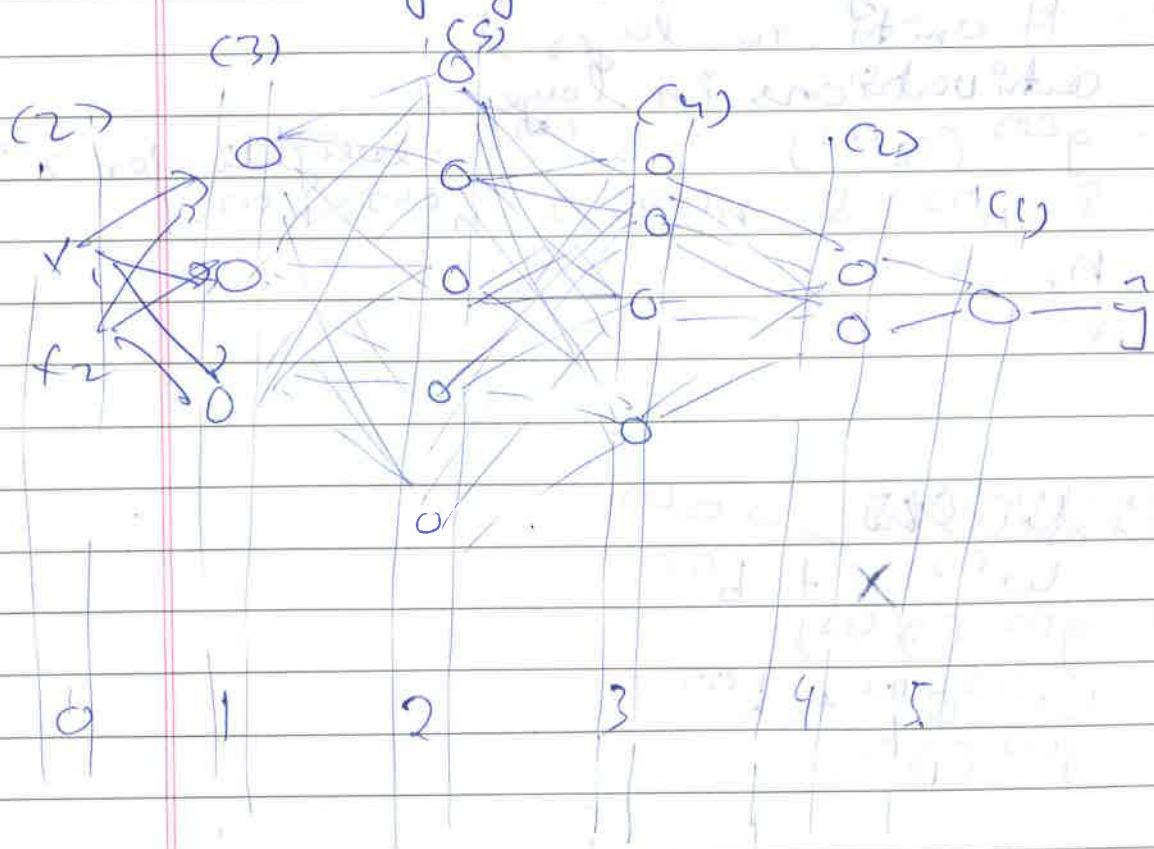
$$A^{(2)} = g^{(2)}(z^{(2)})$$

$$\hat{y} = g(z^{(4)}) = A^{(4)}$$

where

$$z^{(2)} = \begin{bmatrix} z^{(1)(1)} & z^{(1)(2)} & z^{(1)(3)} & z^{(1)(4)} \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

→ Getting your matrix dimensions right.



$$w^{(1)} : (n^{(0)}, n^{(1)}) \quad (5, 2)$$

$$w^{(2)} : (n^{(1)}, n^{(2)}) \quad (2, 3)$$

$$w^{(3)} : (n^{(2)}, n^{(3)}) \quad (3, 1)$$

$$w^{(4)} : (n^{(3)}, n^{(4)}) \quad (1, 1)$$

$$w^{(5)} : (n^{(4)}, n^{(5)}) \quad (1, 1)$$

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



and $\Sigma^{(l)} = [n^{(l)}, 1]$

for authorized versions:-

$$(i) \Sigma^{(l)} = \left[\begin{array}{c} b^{(1)} \\ b^{(2)} \\ b^{(3)} \\ b^{(4)} \\ b^{(5)} \end{array} \right] \quad \left[\begin{array}{c} b^{(1)} \\ b^{(2)} \\ b^{(3)} \\ b^{(4)} \\ b^{(5)} \end{array} \right]$$

$$\Sigma^{(l)} = (n^{(l)}, m)$$

$$(ii) b^{(l)} = \left[\begin{array}{c} b^{(l-1)1} \\ b^{(l-1)2} \\ \vdots \\ b^{(l-1)m} \end{array} \right]$$

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

$$\text{thus } \Delta b^{(l)} = (n^{(l)}, m)$$

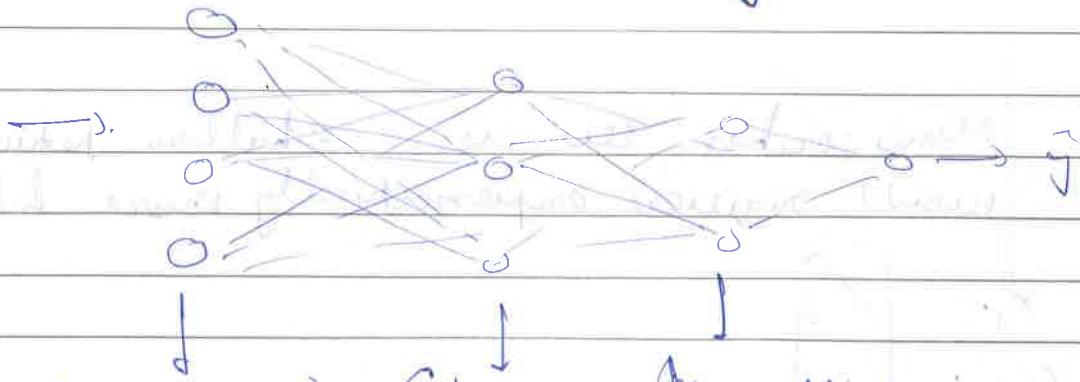
number of nodes
in previous layer

$$(iii) w^{(l)} = \left[\begin{array}{c} \text{node } 1 \\ \text{node } 2 \\ \vdots \\ \text{node } n \end{array} \right] \quad \left\} \begin{array}{l} \text{number of nodes} \\ \text{in previous layer} \end{array} \right\}$$

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

→ Intuition about deep learning.

Picture
of a few
person.



first layer
(may detect the edges)
by putting bunch of edges we may get features like nose eyes
by putting together features we could get a face.

So basically deep learning is deep neural networks

卷之三

Circuit theory and deep learning :-

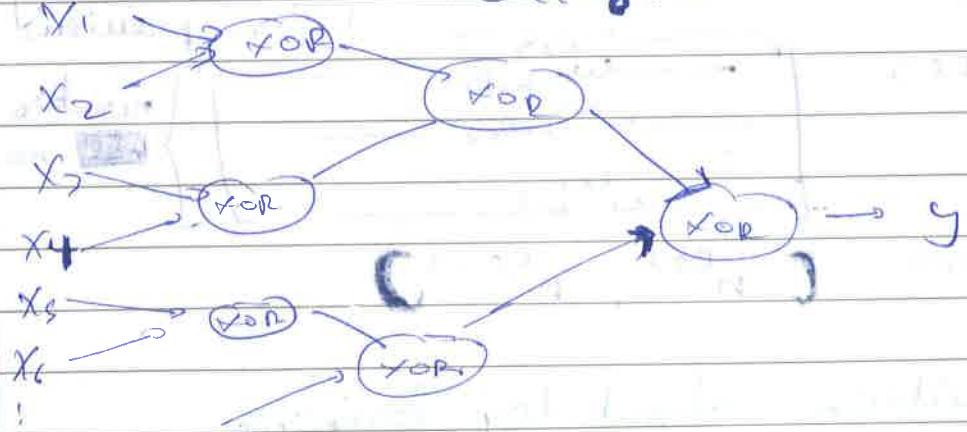
Informally :- There are ~~more~~ functions you can compute with a "small" 1-layer deep neural network than 2-hidden networks require exponentially more hidden units to compute.

59

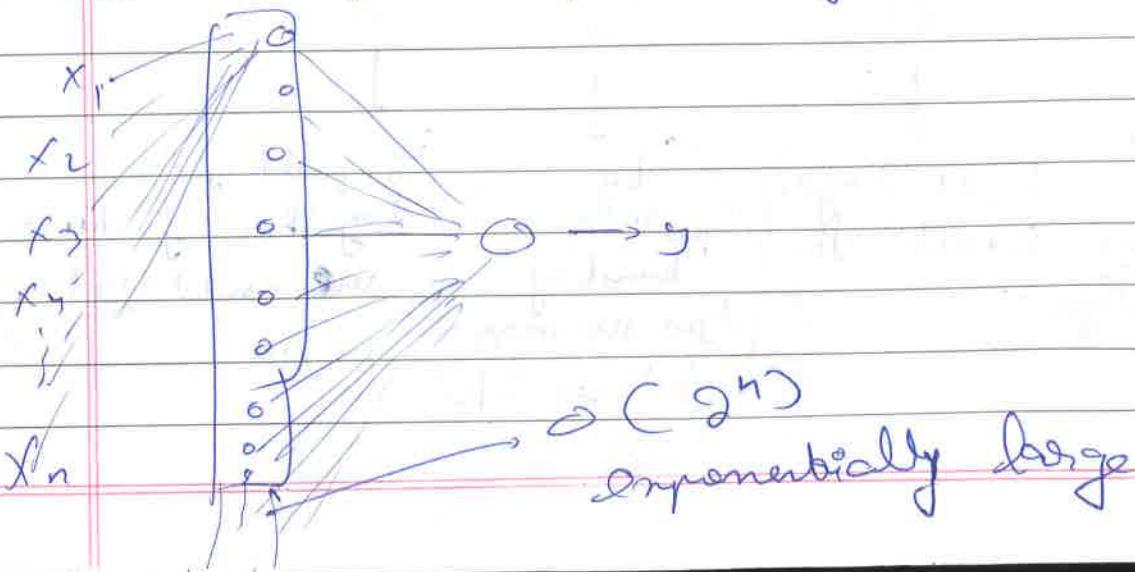
This is when we are using e-gated deep reward networks

$$y = (x_1) \text{XOR} (x_2) \text{XOR} (x_3) \text{XOR} (x_4) \dots \text{XOR} (x_n)$$

$\leftarrow O(\log n) \rightarrow$

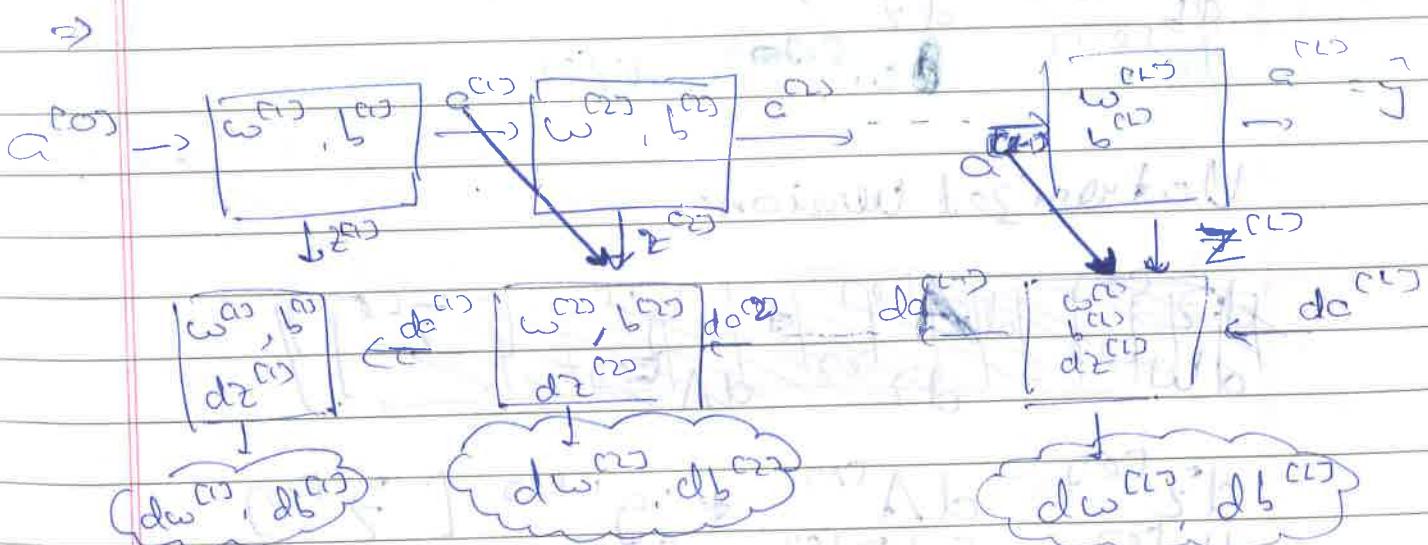
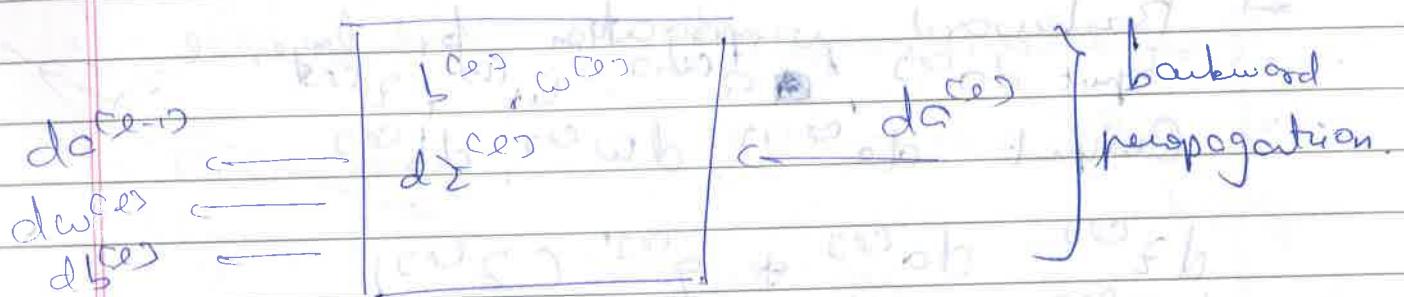
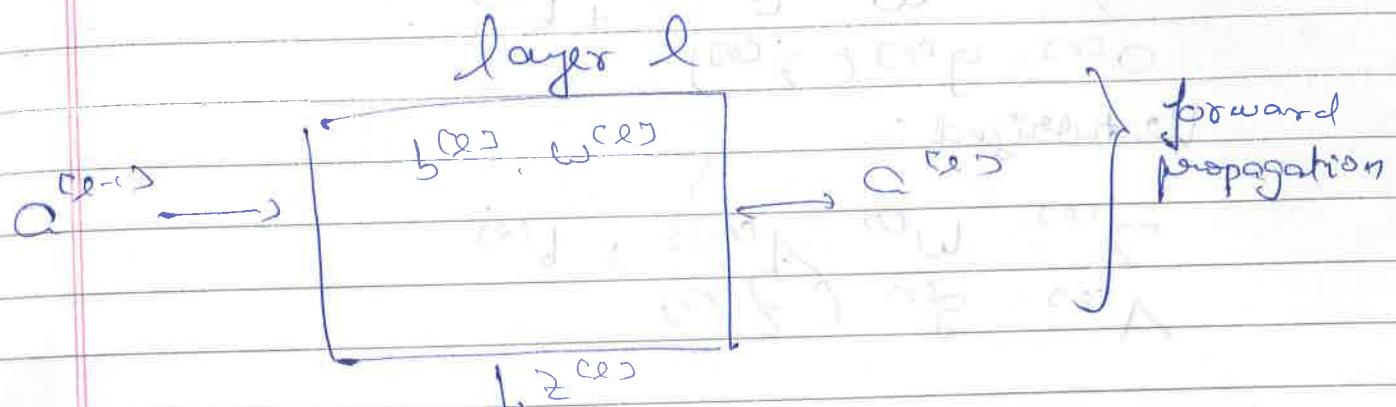


When entering the case of shallow networks we would require exponentially more hidden nodes.



→ Building blocks of deep ~~neural~~ neural network.

→ Forward and backward functions:



$$\begin{aligned} w^{(l)} &:= w^{(l)} - \alpha \cdot dw^{(l)} \\ b^{(l)} &:= b^{(l)} - \alpha \cdot db^{(l)} \end{aligned}$$

→ Forward propagation for layer l .

→ Input $a^{(l-1)}$

→ Output $a^{(l)}$, cache ($z^{(l)}$)

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$

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

Vectorized :-

$$z^{(l)} = w^{(l)} A^{(l-1)} + b^{(l)}$$

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

→ Backward propagation for layer l .

→ Input $da^{(l)}$, $a^{(l-1)}$, $w^{(l)}$, $z^{(l)}$

→ Output $da^{(l-1)}$, $dw^{(l)}$, $db^{(l)}$

$$dz^{(l)} = da^{(l)} * g'(z^{(l)})$$

$$dw^{(l)} = dz^{(l)} \cdot A^{(l-1)T}$$

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

$$da^{(l-1)} = w^{(l)T} \cdot dz^{(l)}$$

Vectorized Version:-

$$\begin{aligned} dz^{(l)} &= da^{(l)} * g'(z^{(l)}) \\ dw^{(l)} &= dz^{(l)} \cdot A^{(l-1)T} \\ db^{(l)} &= dz^{(l)} \end{aligned}$$

$$dz^{(l)} = dA^{(l)} * g'(z^{(l)})$$

$$dw^{(l)} = (dz^{(l)} \cdot A^{(l-1)T}) (1/m)$$

$$db^{(l)} = \frac{1}{m} \cdot \text{np_sum}(dz^{(l)}), \text{axis} = 1, \text{keepdims} = \text{True}$$

$$dA^{(l-1)} = w^{(l)T} \cdot dz^{(l)}$$

→ To calculate

$$dA^{(C2)} = \left(-\frac{y^{(C2)}}{a^{(C2)}} + \frac{(1-y^{(C2)})}{1-a^{(C2)}} \right) \frac{-y^{(C2)}}{a^{(C2)}}$$

→ To calculate

$$da^{(C2)} = \left(-\frac{y^{(C2)}}{a^{(C2)}} + \frac{(1-y^{(C2)})}{1-a^{(C2)}} \right)$$

→ Parameters Vs Hyperparameters.

→ Parameters:- $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}$.

Hyperparameters :- learning rate α , number of iterations

Hidden layer l

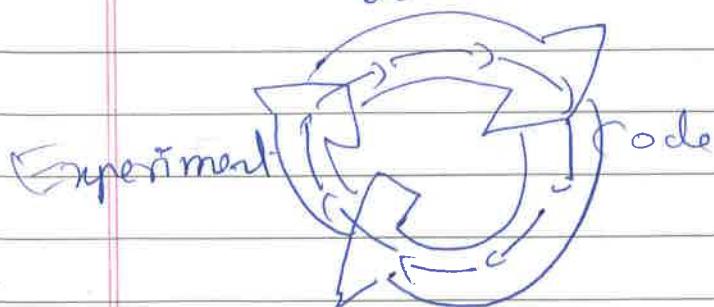
Hidden units $n^{(1)}, n^{(2)}$.

Choice of activation function.

→ Hyperparameters determine the values of parameters.

→ Applied deep learning is a very empirical process.

Idea



It is not very accurate to say that deep learning is like human brain. We know very less about neurons & human brain to make this analogy.

- "Cache" is used to pass variables computed during forward propagation to the corresponding backward propagation step. It contains useful values for backward propagation to compute derivatives.

Eg. like $z^{(l)}$ which is computed in forward propagation is used to compute

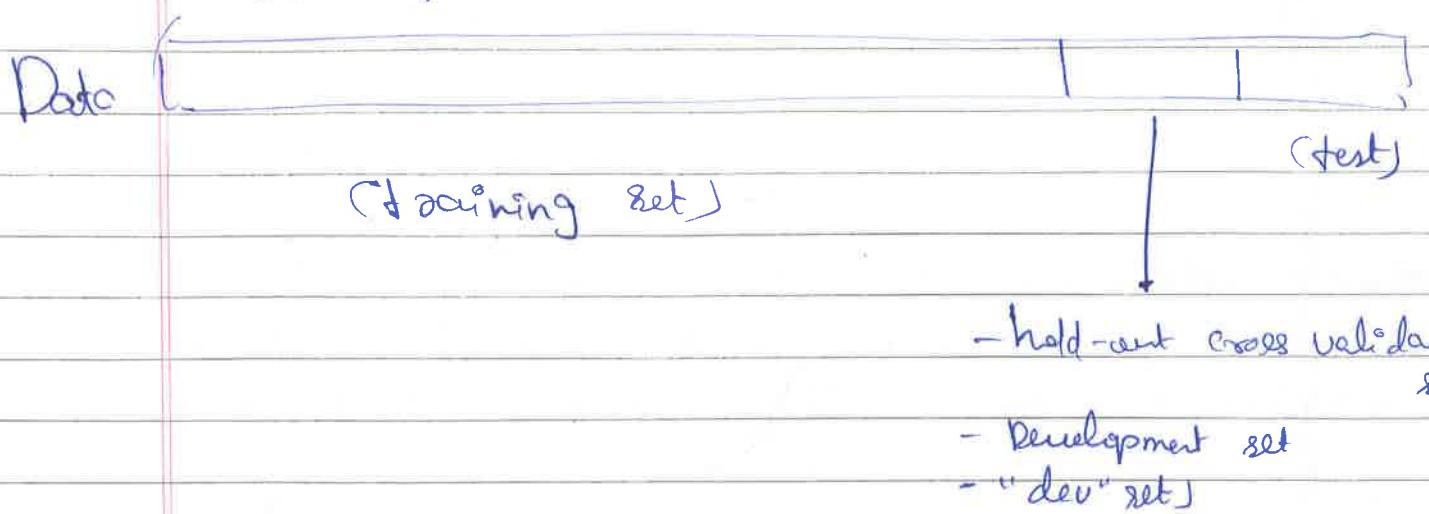
$$d\mathbf{z}^{(l)} = (\mathbf{w}^{(l+1)} \cdot d\mathbf{z}^{(l+1)}) * g'(\mathbf{z}^{(l)})$$

- Some notes from the python code on 1-layer neural network:

Specialization 2 Week 1

PAGE No.	77
DATE	

→ Train | dev | test | sets



⇒ Earlier we used to divide the data as :-
60% / 20% / 20%.

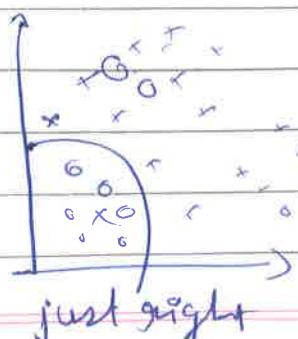
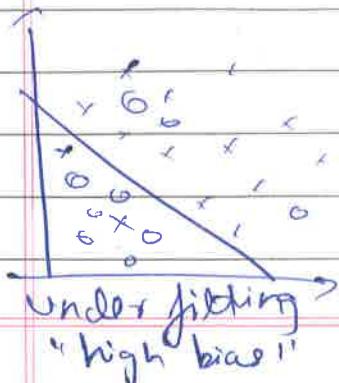
⇒ But now in the age of Big Data, for example when we have 1,000,000 (1 million samples) then even 10,000 samples for dev set/test set would be enough

thus :- 98% / 1% / 1%.

⇒ Make sure that "dev" & "test" set are captured from the same distribution.

⇒ Not having a test set might be okay (only dev set)

→ Bias | variance



Eg Cat classification

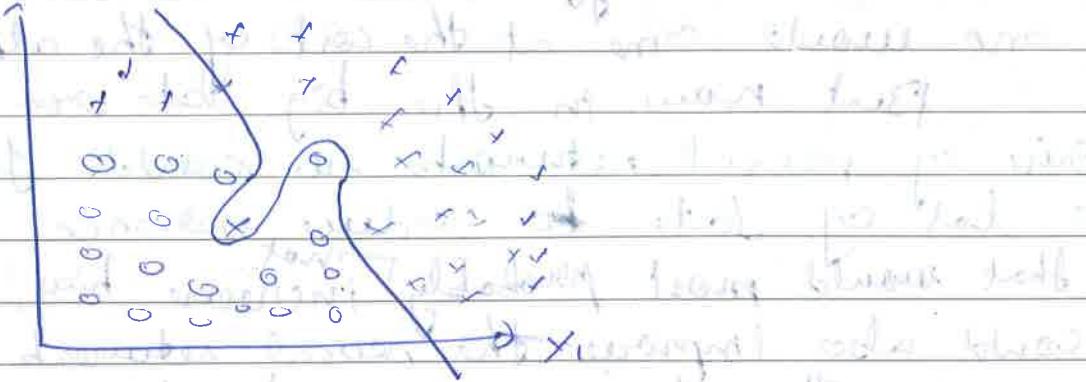
train set error:	17%	15.1%	15%	0.5%
Dev set error:	11%	16%	30%	1%
	(high variance)	(high bias)	(high bias & high variance)	(low bias & low variance)

→ This is assuming that Bayes error ≈ 0.5 & train & Dev set are taken from the same distribution.

→ If Bayes error $\approx 15\%$ (This means that even humans are 15% times wrong while classifying cat image probably because the image is too blur, then → this would not be high bias but seems to be suitable)

→ high bias & high variance

x_1 ~~100 training data points~~



→ Basic recipe for machine learning.

High bias?

(Training data performance)

(No)

(Yes)
high (we don't have high bias)

low (we have high bias)

→ bigger network

→ Train longer (Increase number of iterations)

→ Search for better neural network architecture.

No.

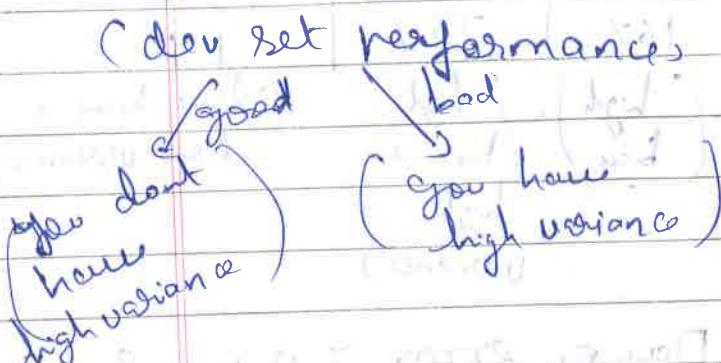
High variance?

Yes.

(More data)

(Regularization)

(get better neural network architecture)



⇒ If you could collect more data then you could reduce variance but not necessarily you could ~~reduce~~ considerably reduce bias.

⇒ Earlier there was concept of bias / variance tradeoff which meant that improving one would come at the cost of the other.

But now in the big data era with rise of neural networks we could feed in a lot of data to reduce variance ~~not~~ that would most probably ^{not} increase bias, we could also improve the neural network architecture. Then the concept of bias / variance tradeoff may not strictly apply.

→ Regularization. (to prevent overfitting)

→ logistic regression:-

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \left((\hat{y}^{(i)}, y^{(i)}) \right) + \frac{1}{2m} \|w\|^2$$

"regularization parameter"

$$\|\mathbf{w}\|_2^2 = \sum_{j=1}^n w_j^2 = \mathbf{w}^T \mathbf{w}$$

\Leftrightarrow \mathbf{L}_2 regularization.

\mathbf{L}_1 regularization:-

$$\text{her per odd } \sum_{i=1}^n |w_i| = \frac{1}{2m} \sum_{i=1}^n |w_i|$$

If we use \mathbf{L}_1 regularization then \mathbf{w} will be sparse, it means that \mathbf{w} would have a lot of zeroes. That would save memory.

\Rightarrow Neural network:-

$$J(\mathbf{w}^{(0)}, \mathbf{b}^{(0)}, \dots, \mathbf{w}^{(L)}, \mathbf{b}^{(L)})$$

$$= \frac{1}{m} \sum_{i=1}^m \ell(g^{(L)}(\mathbf{x}^{(i)}), y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\mathbf{w}^{(l)}\|_F^2$$

$$\|\mathbf{w}^{(l)}\|_F^2 = \sum_{i=1}^m \sum_{j=1}^{n^{(l-1)}} (w_{ij}^{(l)})^2$$

"For benius norm"

dimensions of \mathbf{w} : $(n^{(0)}, n^{(L-1)})$

~~gradient descent~~

~~Newton's method~~

{ p. 700 }

⇒ Earlier we used to do:-

$$dw^{[e]} = (\text{from backpropagation}) + \frac{\lambda}{m} w^{[e]}$$

⇒ Earlier before regularization we used to do:-

$$dw^{[e]} = (\text{the value we got from backpropagation})$$

$$\text{where } dw^{[e]} = \frac{\partial J}{\partial w^{[e]}}$$

now we will do:-

$$dw^{[e]} = (\text{from backpropagation}) + \frac{\lambda}{m} w^{[e]}$$

thus

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

$$\Rightarrow w^{[e]} = w^{[e]} - \alpha (\text{from backpropagation}) + \frac{\lambda}{m} w^{[e]}$$

$$\Rightarrow w^{[e]} = w^{[e]} \left(1 - \frac{\alpha \lambda}{m}\right) - \text{from backpropagation}$$

↳ This is the

This would be less than

thus this is called weighted decay

→ How does regularization prevent overfitting

$$J(w^{[e]}, b^{[e]}) = \frac{1}{m} \sum_{i=1}^m l(y_i, y^{(i)}) + \frac{\lambda}{2m} \sum_{e=1}^E \|w^{[e]}\|_F^2$$

Some

→ now as $\hat{y} \rightarrow 1$, $w^{(c)}$, it is also possible that \boxed{w} of \boxed{b} the values of $w^{(c)}$ ~~would~~ would tend to 0.

Thus then our answers would be more linear.

→ $y \rightarrow 1$ thus $w^{(c)} \downarrow$ then :-

$z^{(c)} = w^{(c)} a^{(c-1)} + b^{(c)}$ would be low.

If we use $g(z) = \tan^{-1}(z)$



as z would be small then we would get $g(z)$ to be more linear. Thus our final ans would tend to be more linear. This is because the ans of every layer would tend to be more linear.

→ $J(\theta) = \sum l(\hat{y}^{(c)}, y^{(c)}) + \frac{\lambda}{2} \sum \|w^{(c)}\|^2$

This would give us
a good
loss

→ Dropout regularization

→ Illustrating dropout on layer $l = 3$

$d_3 = \text{np. random. rand}(\alpha, \text{shape}(a_3, a_3, \text{shape}))$
~~~~~  $\hookrightarrow$  keep\_prob

$a_3 = \text{np. multiply}(a_3, d_3)$

$e_3 = a_3 / \text{keep\_prob}$   $\hookrightarrow$  we do this so that the result of the cost will have the same expected value as without drop-out.

This was "inverted dropout"

→ making predictions at test time.

$$a^{(0)} = x$$

No drop out:

$$z^{(1)} = w^{(1)} a^{(0)} + b^{(1)}$$

$$a^{(1)} = g^{(1)}(z^{(1)})$$

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

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

$$y$$

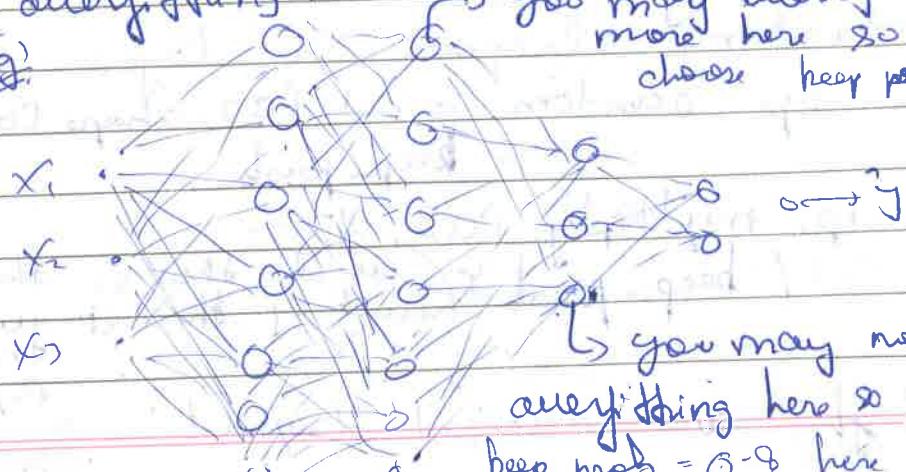
In dropout regularization you randomly drop some nodes for each iteration.

→ Why does drop-out work.

Intuition:- Can't rely on any one feature, so have to spread out weight.

If you are worrying about overfitting ~~on~~ on layer more, then ~~on~~ decrease keep-prob on them and increase keep-prob on those which you ~~think~~ think won't pose a problem for overfitting.

Eg:



→ you may worry about overfitting more here so you may choose keep prob = 0.5 here.

→ you may not worry about overfitting here so you may choose keep prob = 0.9 here.

Here cost function  $J$  is not very much properly defined

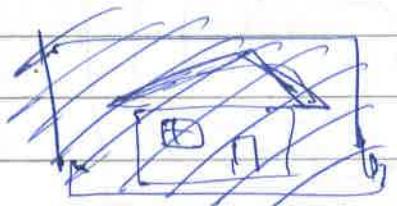
→ Other regularization method.

⇒ Data augmentation.

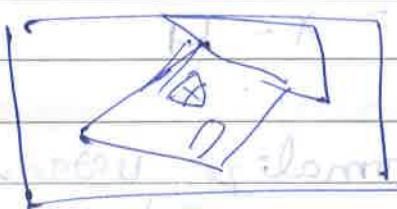
If you have this picture



and you need more data and you don't have then you could do this:-



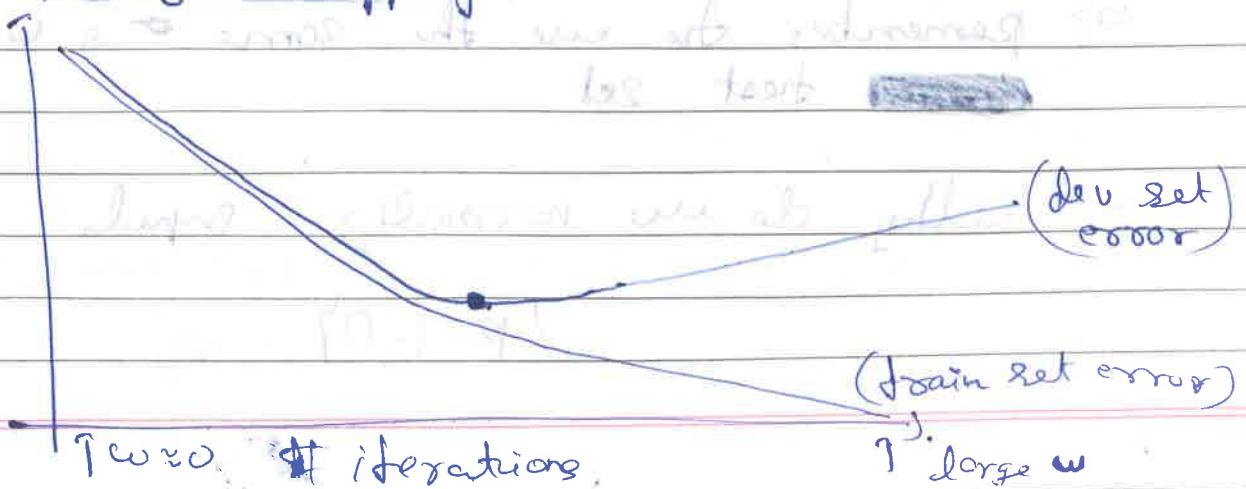
(Clipped sideways)



(Cropped & rotated)

This would not give very much new information but still give a little information and would help a little in regularization.

⇒ Early stopping.



⇒ Orthogonalization.

This means that at one time focus only on one task.

so let's look at three tasks:

(i) Optimize cost function  $J$

(ii) Do not overfit.

So ~~first~~ first task is to concentrate on just decreasing  $J$  & other is to prevent overfitting.

→ Normalizing Inputs

(i) Subtract mean:

$$\bar{u} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x' = x - \bar{u}$$

(ii) normalize variance

$$\text{Variance} = (\text{standard deviation})^2 = \sigma^2$$

$$\Rightarrow \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)})^2 - \bar{u}^2$$

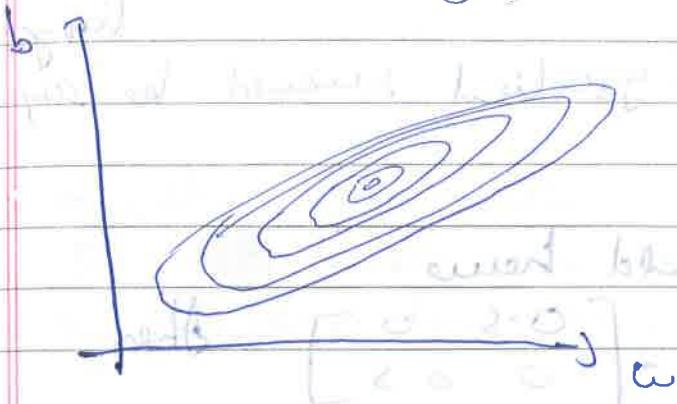
$$x'' = x / \sigma$$

⇒ remember to use the same  $\sigma$  &  $\bar{u}$  on the ~~test~~ test set

→ Why do we normalize inputs:

{p. 7-07}

ii) Un-normalized

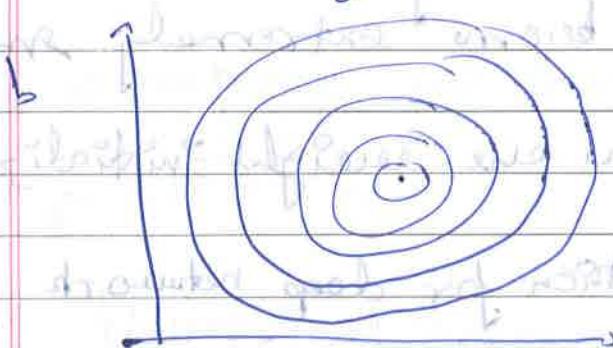


$$x_1: 1 \rightarrow 1000$$

$$x_2: 0 \rightarrow 1$$

graph of  $J$  (cost function)

iii) normalized



$$x_1: -1 \rightarrow 1$$

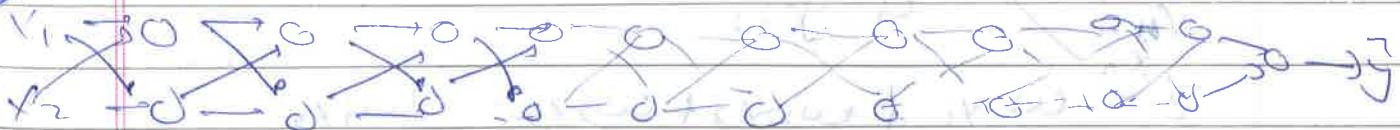
$$x_2: 0 \rightarrow 1$$

graph of  $J$  (cost function)

if we normalize ~~the~~ we would be able to train our model faster.

→ Vanishing / Exploding gradient.

Eg)



$$y = w^{(1)} w^{(2)} w^{(3)} \dots w^{(L-1)} w^{(L)} x$$

assuming that all  $w^{(i)} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$

$$\text{then } y = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x = 1.5^{L-1} x$$

Now here  $\nabla g$  would be exponentially large,

and even the gradient would be exponentially large.

and if we would have

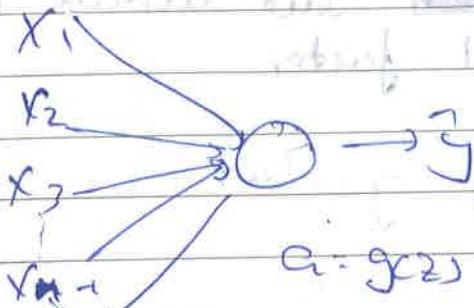
$$w^{(2)} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \text{ then}$$

$\nabla g$  would be extremely small & even gradient would become extremely small

To solve this we use weight initialization.

→ weight initialization for deep network

→ single neuron example.



$$z = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n$$

now if  $n \rightarrow \infty$  is large then we would want  $w_i$  to be small so that  $z$  doesn't become very large.

thus we sometimes define:

$$\text{Var}(w_i) = \frac{2}{n}$$

$$w^{(2)} = \text{np. random. randn(shape)} \otimes \text{np. sqrt}\left(\frac{2}{n^{ce-0}}\right)$$

we prefer to use

$$\text{Var}(\mathbf{w}_i) = \frac{2}{n} \quad \text{when we have } g^{(L)}(\mathbf{z}) = \text{ReLU}(\mathbf{z})$$

in other cases we could use other variance.

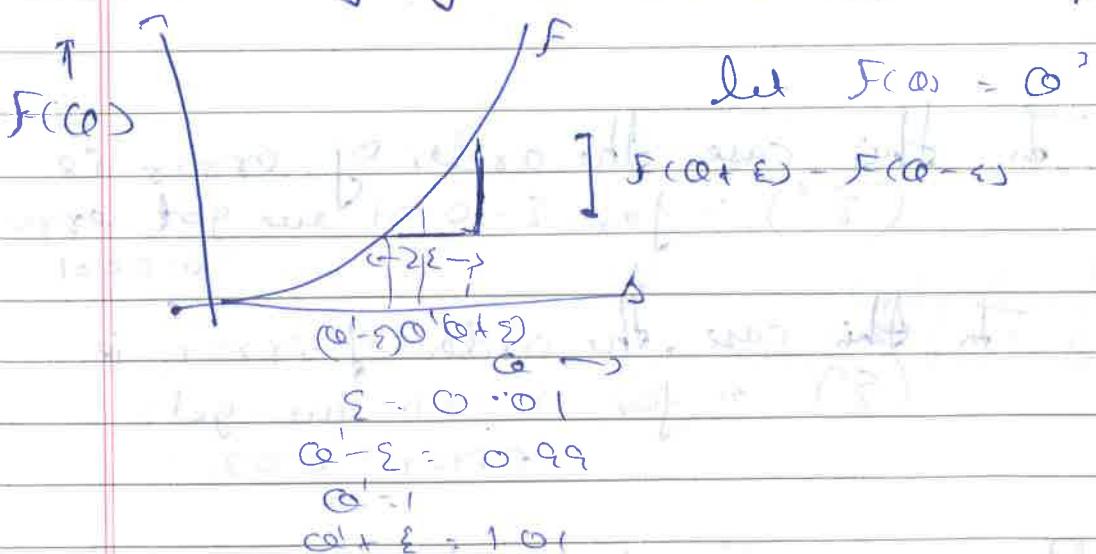
eg. when we have tanh as our activation function then we could use  $\int \frac{1}{n^{(e-1)}}$

This is called Renier Regularization.

sometimes we also use:

$$\sqrt{n^{(e-1)} + n^{(e-2)}}$$

- Numerical Approximation of gradients
- Checking your derivatives computation.



(i)  $F(\theta + \epsilon) - F(\theta - \epsilon) \approx g(\theta)$

Q2

$$\frac{(1.01)^3 - (0.99)^3}{2 \cdot (0.01)} = 3.0001$$

$$g(\theta) = 3\theta^2 = 3$$

approximation error: - 0.0001

(ii) Now trying:

$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}$$

$$= \frac{(1.01)^3 - 1}{0.01} = 3.0301$$

error: 0.02

Then we would prefer to use

$$\left( \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \right) \text{ or gradient over}$$

$$\left( \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \right)$$

→ In this case the order of error is  $(\epsilon^2)$   $\Rightarrow$  for  $\epsilon = 0.01$  we get error = 0.0001

→ In this case the order of error is  $(\epsilon)$   $\Rightarrow$  for  $\epsilon = 0.1$  we get error = 0.03

→ ~~Implementation~~ Gradient checking.

→ Take  $w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}$  and reshape into a big vector  $\theta$ .

$$J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = J(\theta)$$

→ Take  $d\mathbf{w}^{(l)}$ ,  $d\mathbf{b}^{(l)}$ , ...,  $d\mathbf{w}^{(L)}$ ,  $d\mathbf{b}^{(L)}$  and  
reshape into a big vector  $\mathbf{d}\mathbf{o}$ .  
 $\rightarrow (d\mathbf{w}^{(1)}, d\mathbf{b}^{(1)}, \dots, d\mathbf{w}^{(L)}, d\mathbf{b}^{(L)}) = \mathbf{J}(\mathbf{d}\mathbf{o})$

→ Gradient check.  
for each  $i$ :

~~$d\mathbf{o}_{\text{approx}}^{(i)} = \frac{\mathbf{J}(\mathbf{o}_i + \varepsilon, \mathbf{o}_2, \dots, \mathbf{o}_L) - \mathbf{J}(\mathbf{o}_i - \varepsilon, \mathbf{o}_2, \dots, \mathbf{o}_L)}{2\varepsilon}$~~

use  $\varepsilon = 10^{-7}$  most reasonable it's right

check  $\frac{\|d\mathbf{o}_{\text{approx}} - d\mathbf{o}\|_2}{\|d\mathbf{o}_{\text{approx}}\|_2 + \|d\mathbf{o}\|_2}$

$\approx 10^{-7} \rightarrow (\checkmark)$   
 $\approx 10^{-5} \rightarrow (\text{check})$   
 $\approx 10^{-3} \rightarrow (X)$   
 $\downarrow$   
 $\text{Something is wrong.}$

→ Gradient checking implementation notes:-

- Don't use in training - only to debug.
- If algorithm fails grad check, look at component, do try to identify bug.
- Remember regularization.
- Doesn't work with dropout.
- Run w/ random initialization; perhaps again after some time.

~~cost function for L2 regularization~~

→ python code for L2 regularization

let there be a 3 layer neural network with

$$\text{cost} = (\text{old cost}) + (\lambda/2m) * (\text{np.sum(np.square}(w_1)) + \text{np.sum(np.square}(w_2)) + \text{np.sum(np.square}(w_3)))$$

In backpropagation we make some changes only for  $d\omega_1$ ,  $d\omega_2$  &  $d\omega_3$ .

$$d\omega_3 = 1/m (\text{np.dot}(dZ_3, A_2.T)) + (\lambda/m * \omega_3)$$

$$d\omega_2 = 1/m (\text{np.dot}(dZ_2, A_1.T)) + (\lambda/m * \omega_2)$$

$$d\omega_1 = 1/m (\text{np.dot}(dZ_1, X.T)) + (\lambda/m * \omega_1)$$

Search took up the time to look for a dated book for the 1929 edition of the map until a good old ~~good~~ one was found.

work since we're not restricted from doing  
what we want to do and we're not tied  
down to a desk with a computer. I think  
you could do a lot more than you can in  
a classroom.

selected subject areas and can affect

Constitutive model of the soil

turns, looking like a spiral to the eye.

2002. 1. 9. 09:

23)  $\pi$  is open bracket

Oct 1, 1911, 10:11 A.M.

## ANSWER

## week 2 (Specialization 2)

|          |       |
|----------|-------|
| PAGE NO. | 1 / 1 |
| DATE     |       |

→ mini-Batch gradient descent.

→ Batch vs mini batch gradient descent

Vectorization allows you to efficiently compute on  $m$  examples.

→ Suppose make mini-batches from your main training set. E.g. ~~1000~~ You have  $m = 5,000,000$  training examples. So you make mini-batches with one mini-batch having 1000 training examples.

Then we have 5000 mini-batches.

$$X_{(n_x, m)} = [x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(1000)} \ x^{(1001)} \ \dots \ x^{(2000)} \ \dots \ x^{(m)}]$$

$x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(1000)} \ x^{(1001)} \ \dots \ x^{(2000)} \ \dots \ x^{(m)}$

$x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(1000)} \ x^{(1001)} \ \dots \ x^{(2000)} \ \dots \ x^{(m)}$

$$Y_{(n_y, m)} = [y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(1000)} \ y^{(1001)} \ \dots \ y^{(2000)} \ \dots \ y^{(m)}]$$

$y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(1000)} \ y^{(1001)} \ \dots \ y^{(2000)} \ \dots \ y^{(m)}$

$y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(1000)} \ y^{(1001)} \ \dots \ y^{(2000)} \ \dots \ y^{(m)}$

→ Mini-batch gradient descent.

Repeat {

$$\text{for } t = 1, \dots, 5000$$

$$\text{forward prop on } X^{(t)}$$

$$z^{(t)} = w^{(t)} X^{(t)} + b$$

$$A^{(t)} = g^{(t)}(z^{(t)})$$

$$A^{(t)} = g^{(t)}(z^{(t)})$$

$$\text{Compute cost } J = \frac{1}{1000} \sum_{i=1}^l L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_{i=1}^l \|w^{(i)}\|_F^2$$

Backprop to compute gradients  $J^{(t)}$  (using  $x^{(t)}, y^{(t)}$ )

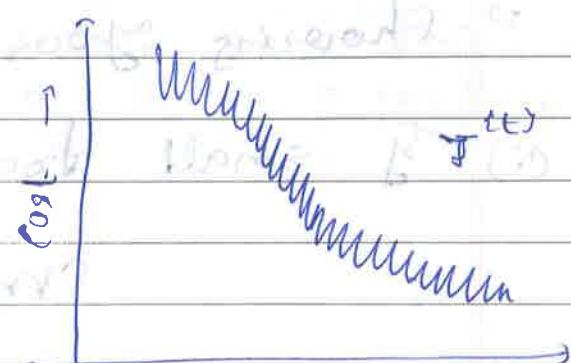
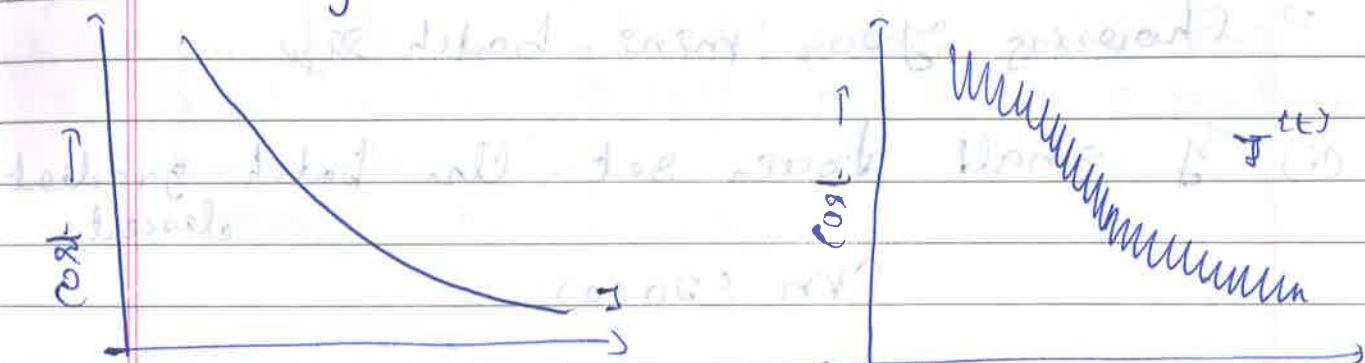
$$w^{(t+1)} = w^{(t)} - \alpha \cdot d w^{(t)},$$

$$b^{(t+1)} = b^{(t)} - \alpha \cdot d b^{(t)}$$

This is 1 epoch "1 pass through training set".

→ Understanding mini-batch gradient descent.

Batch gradient descent.



mini-batch # (b)

Plot  $J^{(t)}$  updated using  $x^{(t)}, y^{(t)}$

⇒ Choosing your mini-batch size.

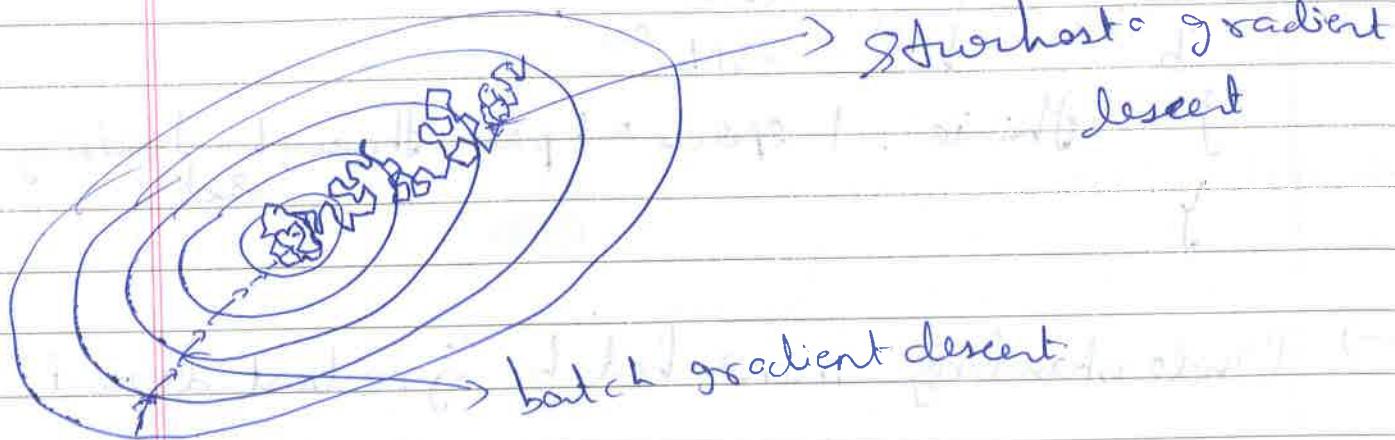
If minibatch size = m: Batch gradient descent  $(x^{(0)}, y^{(0)}) = (y, y)$

If minibatch size = 1: Every example is its own.

Stochastic gradient descent

$$(x^{(0)}, y^{(0)}) = (x^{(0)}, y^{(0)}) - (x^{(0)}, y^{(0)})$$

→ In practice: choose between 1-8m.



⇒ Choosing your mini-batch size.

(i) If small dozen set: Use batch gradient descent.  
( $m \leq 2000$ )

(ii)



Typical mini-batch sizes:

→ 64, 128, 256, 512  
 $2^6$     $2^7$     $2^8$     $2^9$

try to choose the size as a power of 2. This would make the code run faster as the memory is linear in computer ie designed in this way.

(Batch gradient)  
(descent)

→ too long per  
iteration

→ between  
mini-batch size

→ fastest learning  
→ vectorization  
→ makes progress  
without processing  
entire training set

(Stochastic gradient)  
(descent)

→ less error  
speed up from  
vectorization.

→ Exponentially weighted averages

Ex) Temperature in London

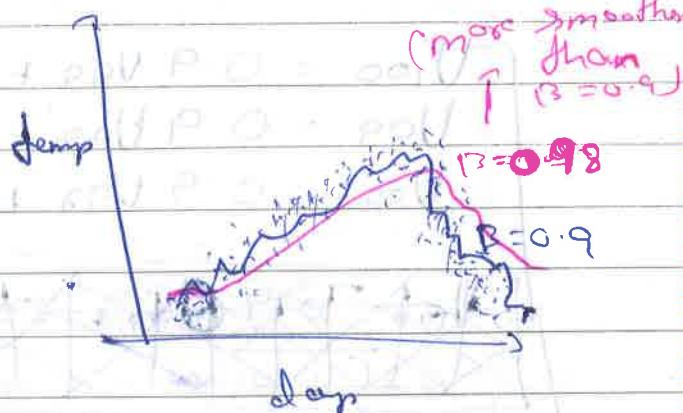
$$O_1 = 40^{\circ} F = 4^{\circ} C$$

$$O_2 = 49^{\circ} F = 9^{\circ} C$$

$$O_3 = 45^{\circ} F = 7^{\circ} C$$

$$O_{180} = 60^{\circ} F = 15^{\circ} C$$

$$O_{181} = 56^{\circ} F = 13^{\circ} C$$



Now let's try to fit a line over this.  
let  $V$  be the value of line at a particular temperature

$$\text{Hence } V = f(\text{days})$$

$$V_0 = 0$$

$$V_1 = (0.9) \cdot V_0 + O_1 \cdot (0.1)$$

$$V_2 = (0.9) \cdot V_1 + O_2 \cdot (0.1)$$

$$V_t = (0.9) \cdot (V_{t-1}) + O_t \cdot (0.1) \rightarrow \text{average}$$

( $V_t$  is approximately  $\text{average } \frac{1}{1-0.9} \text{ days}$ )

Then

$$\text{let } V_t = \beta \cdot (V_{t-1}) + O_t \cdot (1-\beta)$$

( $V_t$  is approximately  $\text{average } \frac{1}{1-\beta} \text{ days}$ )

do if  $\beta = 0.9$

then approximately  $\text{average } 10 \text{ days}$

if  $\beta = 0.98$  then approximately  $\text{average } 50 \text{ days}$

→ Understanding exponentially weighted averages.

$$V_t = \beta \cdot V_{t-1} + (1-\beta) \cdot O_t$$

$$V_{100} = 0.9 V_{99} + (0.1) O_{100}$$

$$V_{99} = 0.9 V_{98} + (0.1) O_{99}$$

$$V_{98} = 0.9 V_{97} + (0.1) O_{98}$$

~~$$\rightarrow V_{100} = 0.9 V_{99} + (0.1) O_{100}$$

$$= 0.9(0.9 V_{98} + (0.1) O_{99}) + (0.1) O_{100}$$

$$= 0.9^2 V_{98} + (0.1) O_{98} + (0.1) O_{100}$$

$$= 0.9^3 V_{97} + (0.1) O_{97} + (0.1) O_{100}$$

$$\vdots$$~~

~~$$V_{100} = 0.9 V_{99} + (0.1) O_{100}$$~~

$$\rightarrow V_{100} = (0.1) O_{100} + 0.9 (0.1) O_{99} + 0.9 (0.1) O_{98} + 0.9 (0.1) O_{97} + \dots$$

$$\Rightarrow V_{100} = (0.1) O_{100} + (0.1)(0.9) O_{99} + (0.1)(0.9)^2 O_{98} + (0.1)(0.9)^3 O_{97} + (0.1)(0.9)^4 O_{96}$$

$$0.9^{10} \approx 0.35 \approx \frac{1}{e}$$

$$(1-\alpha)^{10} \approx \frac{1}{e}$$

→ Implementing exponentially weighted average

$$V_0 = 0$$

$$V_1 = \beta \cdot V_0 + (1-\beta) \cdot O_1$$

$$V_2 = \beta \cdot V_1 + (1-\beta) \cdot O_2$$

$$V_3 = \beta \cdot V_2 + (1-\beta) \cdot O_3$$

hp. T. O]

No = 0

Repeat &

get next  $U_t$

$$U_t = \beta \cdot U_{t-1} + (1-\beta) \cdot \alpha_t$$

Y

→ Bias correction is exponentially corrected averages.

⇒ Bias correction.

$$U_{t-1} = \frac{U_t}{1-\beta}$$

See the thing is that initially we initialize  $U_0 = 0$  so when we calculate  $U_1 = 0.9 U_0 + 0.1 \alpha_1$  then  $U_1$  would be much lesser than the actual value? So what we could do is that we could write  $U_t = \frac{U_t}{1-\beta}$

And this would happen for quite a few initial values.

So for example

$$\beta = 0.9$$

~~$$U_1 = 0.9 U_0 + 0.1 \alpha_1$$~~

~~$$U_1 = \frac{U_1}{1-0.9} = 0.1 U_1 + 0.9 \alpha_1$$~~

~~$$U_1 = 0.9 U_0 + 0.1 \alpha_1$$~~

$$U_1 = \frac{0.9 + 0.1 \alpha_1}{1-0.9} = 0.9 + 0.1 \alpha_1$$

$$\Rightarrow U_1 = 0.9 U_0 + 0.1 \alpha_1 = 0.9 + 0.1 \alpha_1$$

$$\Rightarrow U_2 = 0.9(U_1) + 0.1 \alpha_2 = 0.9(0.9 + 0.1 \alpha_1) + 0.1 \alpha_2$$

$$\Rightarrow U_2 = \frac{0.9(0.9 + 0.1 \alpha_1) + 0.1 \alpha_2}{1-0.9^2} = 0.81 + 0.1 \alpha_2$$

$$\Rightarrow U_2 = \frac{0.9 \cdot 0.9 + 0.1 \alpha_1 + 0.1 \alpha_2}{0.19}$$

let  $\beta = 0.98$

|          |        |
|----------|--------|
| PAGE NO. | 11     |
| DATE     | 1/1/14 |

$$U_0 = 0$$

$$U_1 = (0.98)(U_0) + (0.02)(0_1)$$

$$\frac{U_1}{1 - (0.98)} = \frac{(0.98)(U_0)}{1 - 0.98} + (0.02)(0_1)$$

$$\frac{U_1}{0.02} \rightarrow 0.98(U_0) + (0.02)(0_1)$$

$\rightarrow 49U_0 + 0_1$  → use this value to plot graph

$$\text{now } U_2 = (0.98)(U_1) + (0.02)(0_2)$$

$$\Rightarrow U_2 = (0.98)(0.98)U_0 + (0.02)(0_1) + (0.02)(0_2)$$

$$\Rightarrow U_2 = (0.98)^2 U_0 + (0.98)(0.02)0_1 + (0.02)(0_2)$$

$$\frac{U_2}{1 - (0.98)^2} = \frac{U_2}{0.0396} = \frac{(0.98)^2 U_0}{0.0396} + \frac{(0.98)(0.02)0_1}{0.0396} + \frac{(0.02)(0_2)}{0.0396}$$

$$= 24.282U_0 + 0.49490_1 + 0.50500_2$$

→ plot this on the graph

→ Gradient descent with momentum



Now so to prevent the problem we want a :-  
( $\uparrow$  slower learning rate in this direction)  
( $\leftarrow$  larger learning rate in this direction)

On iteration  $t \rightarrow$

|          |     |
|----------|-----|
| PAGE NO. | 111 |
| DATE     |     |

Compute  $dw$ ,  $db$  on current mini-batch:

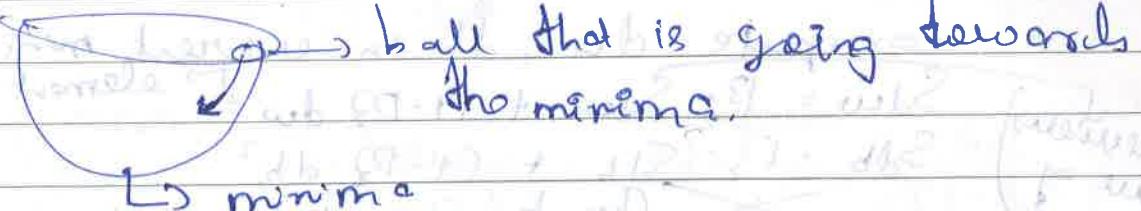
$$Udw = \beta \cdot Udw + (1-\beta) dw$$

$$Udb = \beta \cdot Udb + (1-\beta) \cdot db$$

$$w := w - \alpha \cdot Udw$$

$$b := b - \alpha \cdot Udb$$

→ you can understand this by the analogy of :-



so for this ball we can see the equation as:-

$$Udw = \beta \cdot Udw + (1-\beta) \cdot dw$$

↳ Velocity  $\rightarrow$   $\beta$  acceleration

↳ friction

just to get a better intuition

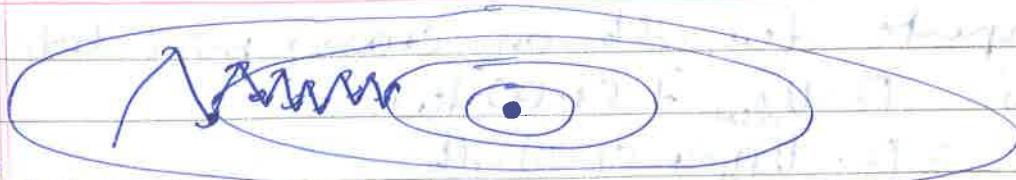
→ We initialize  $Udw = 0$  &  $Udb = 0$  & we have hyperparameters  $\alpha$  &  $\beta$ .

→ Often we don't implement bias correction, because we think that it is not that necessary here

→ RMS prop.

{P.T.O}

we want to reduce the disturbance here



- we want to slow down learning in this direction  
← we want to speed up learning in this direction

On iteration  $t$ :

compute  $d\omega$ ,  $db$  on current mini-batch

$$Sd\omega = \beta_2 Sd\omega + (1-\beta_2) d\omega^2 \quad \text{element wise} \rightarrow \text{Global}$$

$$Sdb = \beta_2 Sdb + (1-\beta_2) d b^2 \quad \rightarrow \text{Relative}$$

Then  $b$  means  $\rightarrow$  the direction of  $b$  in  $\omega$   $\rightarrow$  joining  $b$  coefficients  $\rightarrow$  learning

$$\omega = \omega - \alpha \cdot d\omega \quad b = b - \alpha \cdot db$$

(Previous  
value of  
 $Sd\omega$ )

(as this is relatively  
small, updates of  $\omega$  would  
be relatively larger)

$$\sqrt{Sd\omega + \epsilon}$$

(as this is  
relatively large  
updates of  $b$  would be  
relatively smaller)

→ we add this small  
number to ensure that the denominator  
doesn't become 0 ever.

→ Adam optimization algorithm

Initialization:-

$$Vd\omega = 0, Sd\omega = 0, Vdb = 0, Sdb = 0$$

(P-7.0)

On iteration  $t$ :

Compute  $d\omega$ ,  $db$  using current mini-batch:

$$U_{dw} = \beta_1 \cdot U_{dw} + (1 - \beta_1) d\omega \quad \text{momentum}$$

~~$$U_{db} = \beta_1 \cdot U_{db} + (1 - \beta_1) db$$~~

$$S_{dw} = \beta_2 \cdot S_{dw} + (1 - \beta_2) d\omega^2 \quad \text{"RMSprop"}$$

$$S_{db} = \beta_2 \cdot S_{db} + (1 - \beta_2) db^2$$

$$U_{dw}^{\text{corrected}} = \frac{U_{dw}}{(1 - \beta_1^t)}, \quad U_{db}^{\text{corrected}} = \frac{U_{db}}{(1 - \beta_1^t)}$$

$$S_{dw}^{\text{corrected}} = \frac{S_{dw}}{(1 - \beta_2^t)}, \quad S_{db}^{\text{corrected}} = \frac{S_{db}}{(1 - \beta_2^t)}$$

$$\omega = \omega - \alpha \cdot \frac{U_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}}} + \epsilon$$

$$b = b - \alpha \cdot \frac{U_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}}} + \epsilon$$

$\Rightarrow$  Hyperparameters choice :-

$\alpha$ : needs to be tuned.

$$\beta_1 : 0.9 \quad (\text{d}\omega)$$

$$\beta_2 : 0.999 \quad (\text{d}\omega^2)$$

$$\epsilon : 10^{-8}$$

$\begin{cases} A & P \\ \hookrightarrow \text{Adaptive} & \hookrightarrow \text{Moment} \end{cases} \uparrow \text{M}$  estimation.

## → Learning rate decay.

If we gradually reduce learning rate alpha then in the initial stages when our learning rate is still high our learning will be high but as we reach our minima we could reduce our learning rate so that our J value would just not keep boggling around minima like:-



=> 1 epoch :- 1 iteration / 1 pass through the data.

$$\alpha = \frac{1}{1 + (\text{decay rate}) * (\text{epoch number})}$$

so if  $\alpha_0 = 2$  & decay rate = 1  
then:-

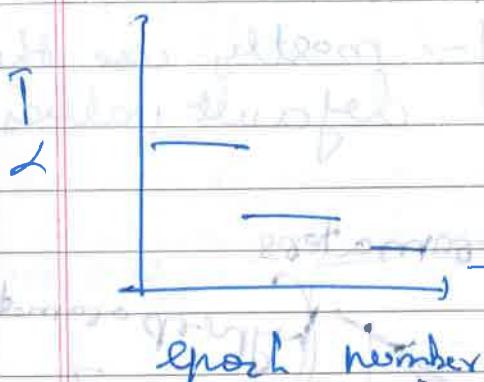
| Epoch no. | $\alpha$ |
|-----------|----------|
| 1         | 0.1      |
| 2         | 0.067    |
| 3         | 0.05     |
| 4         | 0.04     |
| ...       | ...      |

⇒ Other learning rate decay methods :-

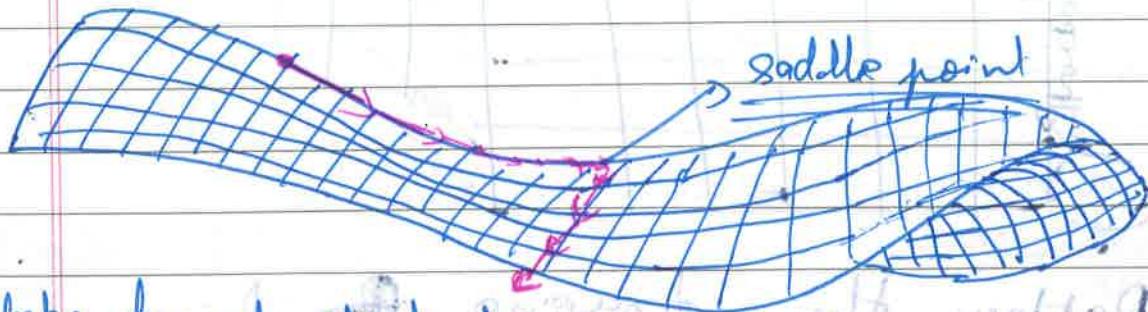
$$\alpha = (0.95)^{\text{epoch-number}} \cdot (\alpha_0)$$

↳ Exponential decay.

$$\alpha = \frac{k \cdot \alpha_0}{\text{epoch-number}}$$



→ The problem of local optima in neural network.



Unlike to get stuck in a bad local optima. This is because there are so many dimensions that it would ultimately go down. It is going down after being stuck at a saddle point for some time. Put these saddle points in plateau can make learning very slow.

Assume that this is a graph of J value over various parameters.

## Specialization 2 (Cont'd)

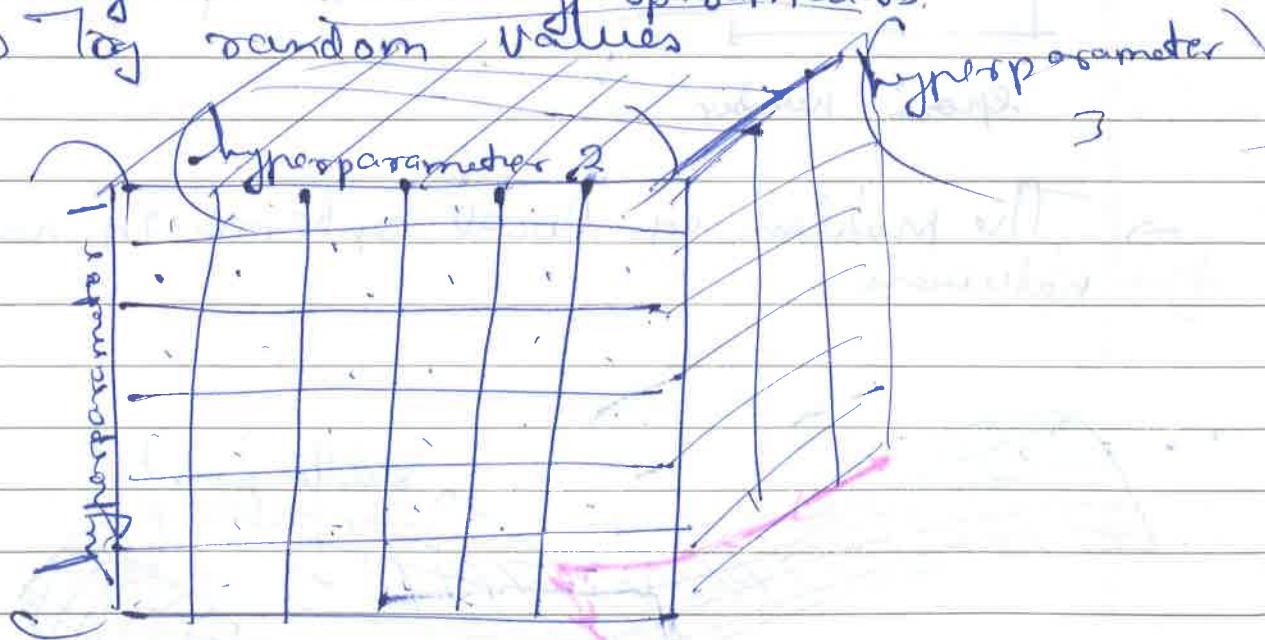
|          |    |
|----------|----|
| PAGE No. | 11 |
| DATE     |    |

### → Tuning hyperparameters

- $\alpha$  → most imp to tune.
  - $\beta$  (for momentum)
  - # hidden units
  - mini-batch size
  - H layers
  - learning rate decay
  - $\beta_1, \beta_2, \epsilon$
  - 0.9
  - $10^{-8}$
  - 0.99
- } Second most important to tune.
- } Third most imp to tune.
- } mostly use these default values.

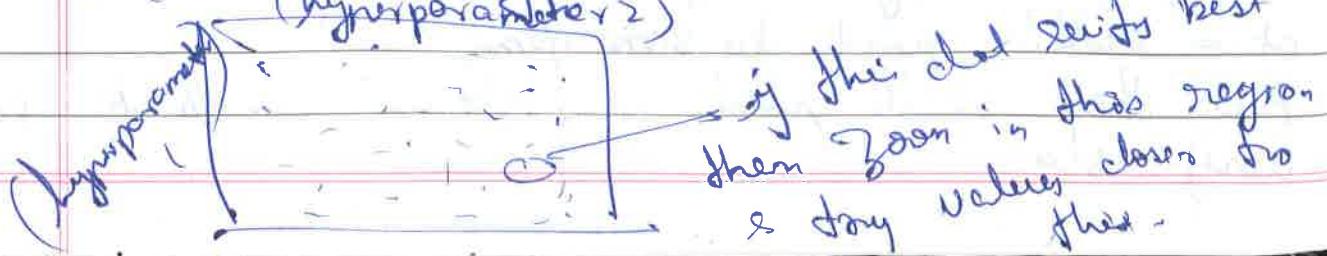
### ⇒ Ways to tune hyperparameters.

#### ⇒ Try random values



Rather than trying ~~the~~ hyperparameters Values in ascending Values systematically try random values of hyperparameters.

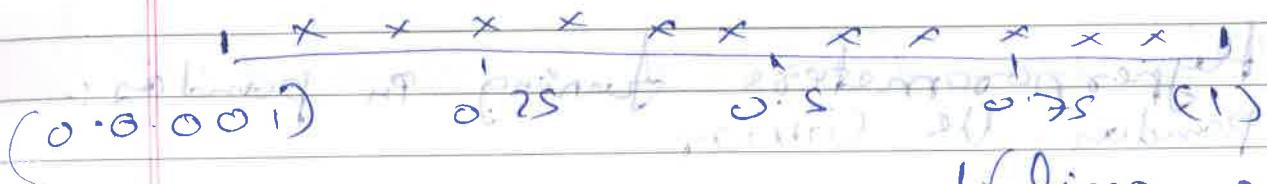
#### ⇒ Coarse to fine



→ Using an appropriate scale to pick hyperparameters.

Eg.)  $\alpha = 0.0001$  to

instead of searching on a linear scale between  $0.0001$  to  $1$  we should prefer to search  $\alpha$  on a logarithmic scale.



(linear scale).

0.0001 0.001 0.01 0.1 1

$\alpha = 10^{\text{rand}}$

for this we write as:

$\alpha = 10^{-4} * \text{np.random.rand}()$

$\alpha = 10^{-4}$

Eg.)  $B$  used to calculate exponentially weighted averages.

$$B = 0.9 - \text{np.random}() - 0.999$$

Here also we don't want to search on a linear scale.

So we will use logarithmic scale on  $(1-B)$

This is because results are very sensitive to change in the value of  $B$  when it is close to one.

$$1 - 10^{-2} = 0.1 \text{ and } 10^{-3} = 0.001$$

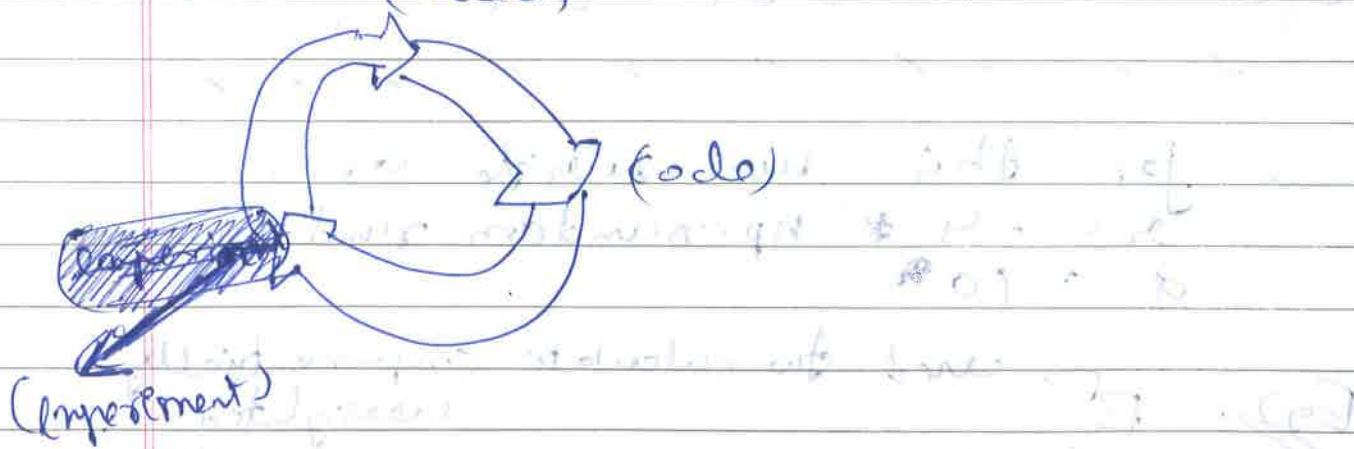
$$1 - 10^{-2} = 0.1 \text{ and } 10^{-3} = 0.001$$

$10^{-2}$  is no good for  $10^{-3}$

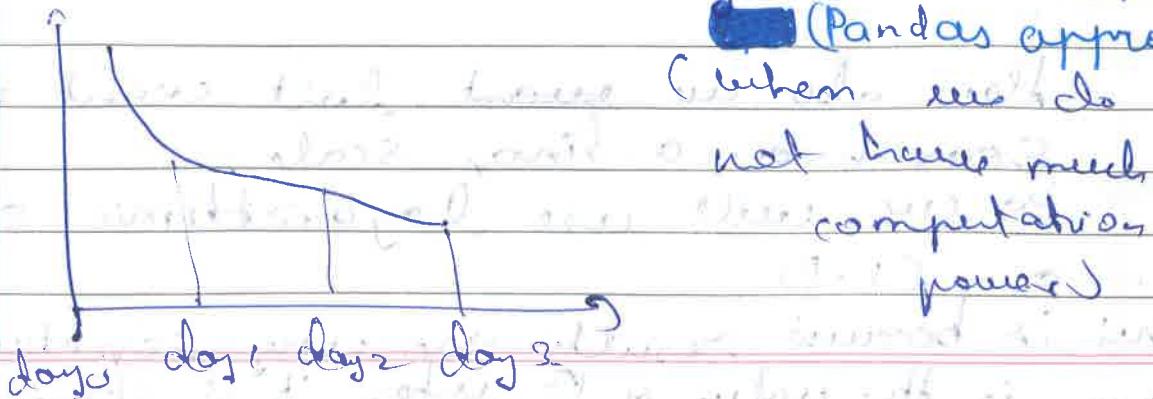
Because our first number measured with  $10^{-2}$  is  $10^{-3}$  and after  $10^{-13} = 10^{-2}$   
 $10^{-3} = 1 - 10^{-2}$

→ Hyperparameters tuning on pandas vs Pandas vs Caviar.

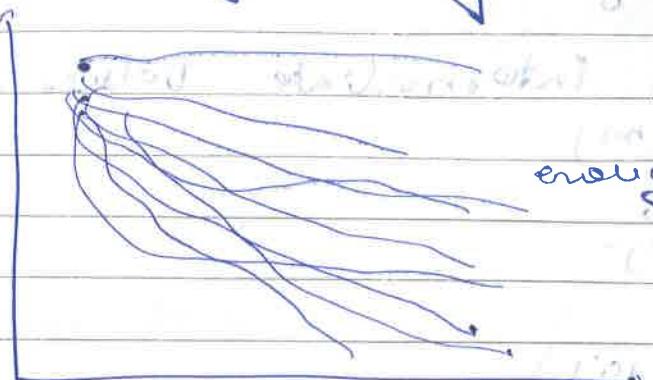
→ Re-test hyperparameters on occasionally (Idea)



→ Either we could babysit one model (Pandas approach)



- Training many models in parallel. (Master approach)



(When we have  
enough computers  
then we do  
this)

- Batch normalization.

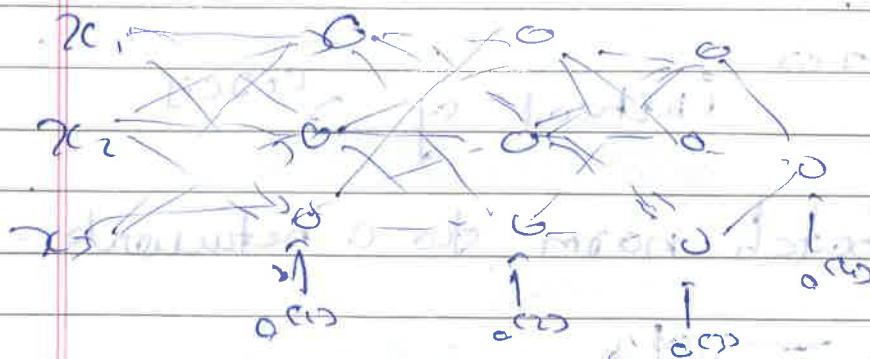
- Normalizing input to speed up learning.

$$\bar{x} = \frac{1}{m} \sum x^{(i)}$$

$$s^2 = \frac{1}{m} \sum (x^{(i)} - \bar{x})^2$$

$$x = \frac{x - \bar{x}}{s} \rightarrow \text{After doing this we have}$$

time on by  $\bar{x}$  and standard deviation  $s$ !



If we normalize  $x_1, x_2, x_3$  then it would become faster to calculate  $w_1, w_2, w_3$ . Similarly if we could normalize  $x^{(2)}$  then it would become faster to calculate  $w_2, b_2$ . This is what batch norm does.

⇒ Implementing batch norm.

Given some intermediate values in N.N

$$z^{(c)} = \underbrace{z^{(c)}_{1,1}, \dots, z^{(c)}_{1,m}}_{\text{before scale}} \quad \underbrace{z^{(c)}_{2,1}, \dots, z^{(c)}_{2,m}}_{\text{after scale}}$$

$$\mu = \frac{1}{m} \sum z^{(c)ij}$$

$$\sigma^2 = \frac{1}{m} \sum (z^{(c)ij} - \mu)^2$$

$$z^{(c)norm} = \frac{z^{(c)ij} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

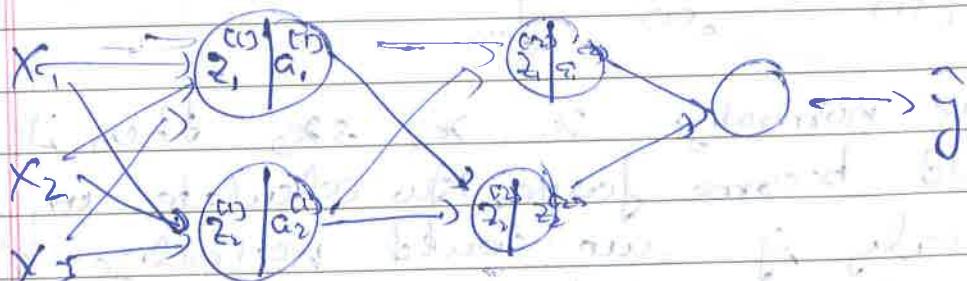
6 we add this so that denom. never becomes 0.

$$\tilde{z}^{(c)ij} = \gamma z^{(c)ij}_{norm} + \beta$$

where  $\gamma$  &  $\beta$  are learning parameters of the model.

use  $\tilde{z}^{(c)ij}$  instead of  $z^{(c)ij}$

→ Adding batch norm to a network :-



$$X \xrightarrow{w^{(l)}, b^{(l)}} Z \xrightarrow{B^{(l)}, \gamma^{(l)}} \tilde{Z} \xrightarrow{C^{(l)}, g^{(l)}} (Z^{(l)})$$

red shaded area Batch norm (BN)

$$X \rightarrow \text{mean} \text{ and } \text{std dev} \rightarrow \tilde{Z} = \frac{Z - \text{mean}}{\text{std dev}}$$

$$\text{mean} = \frac{1}{n} \sum Z^{(l)} \quad \text{std dev} = \sqrt{\frac{1}{n} \sum (Z^{(l)} - \text{mean})^2}$$

Parameters :-  $w^{(l)}, b^{(l)}, w^{(l)}, b^{(l)}, \dots, w^{(l)}, b^{(l)}$

$$B^{(l)}, \gamma^{(l)}, \{B^{(l)}, \gamma^{(l)}\} = \{B^{(l)}, \gamma^{(l)}\}$$

$\text{mean} = \text{mean} \quad \text{std dev} = \text{std dev}$

$$B^{(l)} = \gamma^{(l)} - d \cdot d \gamma^{(l)}$$

⇒ Working with mini-batches

$$X \xrightarrow{w^{(l)}, b^{(l)}} Z \xrightarrow{B^{(l)}, \gamma^{(l)}} \tilde{Z} \xrightarrow{C^{(l)}} (Z^{(l)}) = a^{(l)}$$

$X^{(2)}, X^{(3)}, \dots$  ~~total width no labels need are~~

~~After normalizing batch~~

$$Z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$

After cancel this term. This is because anyway we are going to subtract mean so this will cancel out.

$$Z^{(l)} \text{ normalized } Z^{(l)} - \bar{Z}^{(l)}$$

$$\text{std deviation } \sqrt{\sigma^2 + \epsilon} \text{ adding with mean}$$

$$\tilde{Z}^{(l)} = \gamma^{(l)} Z^{(l)} + B^{(l)}$$

so we are going to do this in parallel with batch size

→ Implementing gradient descent  
for  $t=1 \dots$  number of mini-batches

Compute forward prop on  $X^{(t)}$

In each hidden layer, use Batch norms  
replace  $z^{(c)}$  with  $\tilde{z}^{(c)}$ , use backprop  
to compute  $d\omega^{(c)}$ ,  $d\beta^{(c)}$ ,  $d\gamma^{(c)}$ ,  $dc^{(c)}$   
(not needed)

Update parameters:-

$$\omega^{(c)} := \omega^{(c)} - \alpha \cdot d\omega^{(c)}$$

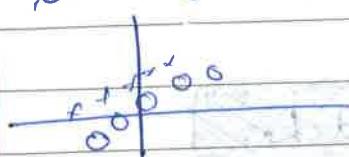
$$\beta^{(c)} := \beta^{(c)} - \alpha \cdot d\beta^{(c)}$$

$$\gamma^{(c)} := \gamma^{(c)} - \alpha \cdot d\gamma^{(c)}$$

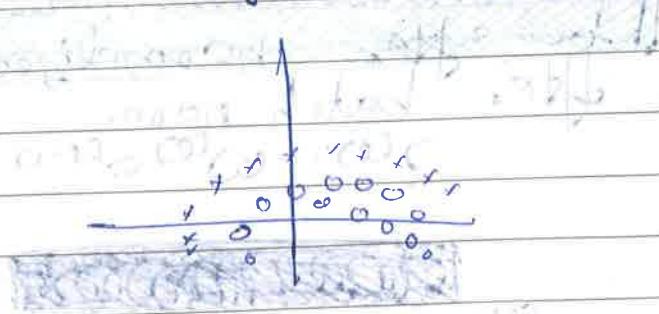
This works with momentum, RMS prop, Adam.

→ Why does batch norm  work.

Ex) we train a model on these data:



and we try our model on data:



so our model may not be able to fit so well. This problem is called covariate shift.

↳ this means

that if we train a model on  $X \rightarrow Y$  & we change  
the distribution of  $X$  then we have to  retrain  
our model.

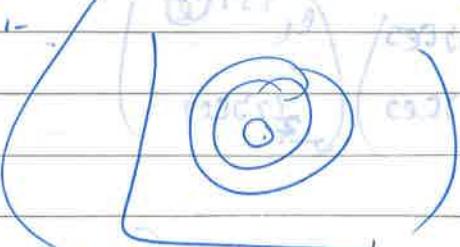
The basic idea behind batch normalization is to limit covariate shift by normalizing the activations of each layer (transforming the inputs to be mean 0 and unit variance). This, supposedly, allows each layer to learn on a more stable distributions of inputs, and would thus accelerate the training of the network.

Batch normalization also induces regularization but less rarely use batch normalization for the purpose of regularization.

→ Actually some recent papers state that limiting the covariate shift is not the primary reason, probably it isn't even an important reason why batch normalization works.

You can I think safely say that batch normalization must be because it does this:

→ See this:



Then making gradient descent & learning faster.

→ Batch norm at test time.

$$\bar{u} = \frac{1}{m} \sum_i z^{(i)}$$

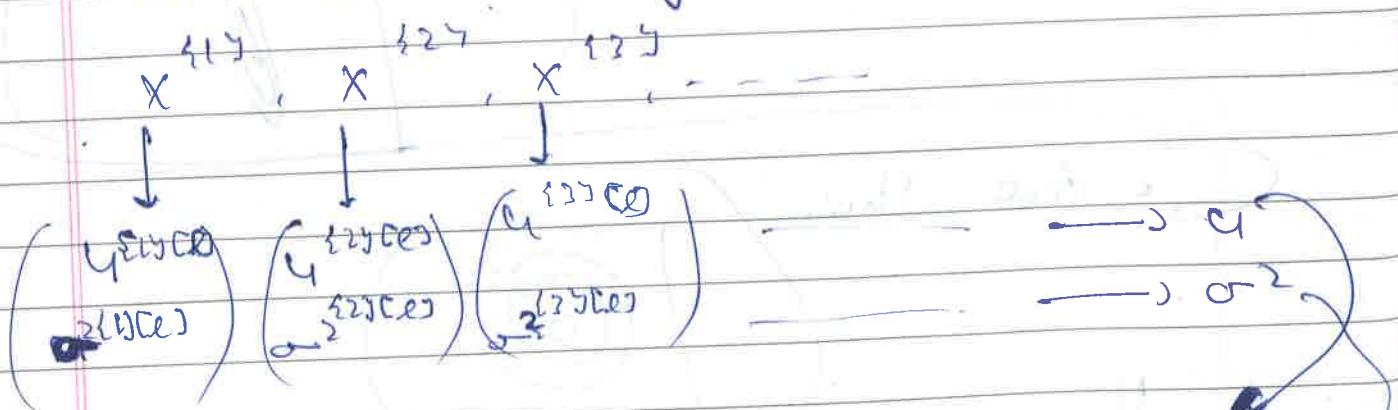
$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \bar{u})^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \bar{u}}{\sqrt{\sigma^2 + \epsilon}}$$

$$z^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

Now during training we can calculate  $\bar{u}$  &  $\sigma^2$  as we have a lot of examples or some examples <sup>while</sup> training. But this might not be case while testing.

Then we do this while testing :-  
estimate  $\bar{u}$ ,  $\sigma^2$  using exponentially weighted average (across mini batches)  
store the values of  $\bar{u}$  &  $\sigma^2$  for all batches



(exponentially weighted averages of stored  $\bar{u}$ )

(exponentially weighted averages of stored  $\sigma^2$ )

And then calculate

$$z_{\text{norm}} = \frac{z - \bar{z}}{\sqrt{\sigma^2 + \epsilon}}$$

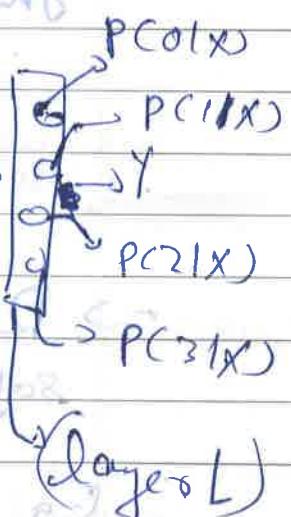
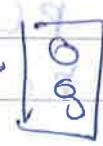
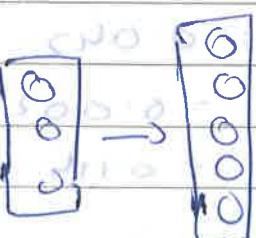
$$\tilde{z} = \gamma z_{\text{norm}} + \beta$$

→ Softmax regression

(Generalization of logistic regression so that we can classify multiple classes instead of just binary classes)

→ 0 → none of the above  
1 → cat  
2 → dog

3 → check



$$z^{(l)} = \omega^{(l)} a^{(l-1)} + b^{(l)} \quad (4.1)$$

Activation function:

$$t = e^{z^{(l)}}$$

$$a^{(l)} = \frac{e^{z^{(l)}}}{\sum_{i=1}^q t_i}$$

$$a^{(l)} = \frac{t_i}{\sum_{i=1}^q t_i}$$

$$\text{Ex2} \quad Z^{(0)} = \begin{bmatrix} 8 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$$t = \begin{bmatrix} e^s \\ e^{2s} \\ e^{-s} \\ e^s \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$a^{(1)} = \frac{t}{176.3}$$

$$\text{Thus } p(0|X) = e^s / (e^s + e^{-s}) = 0.842 = 84.2\%$$

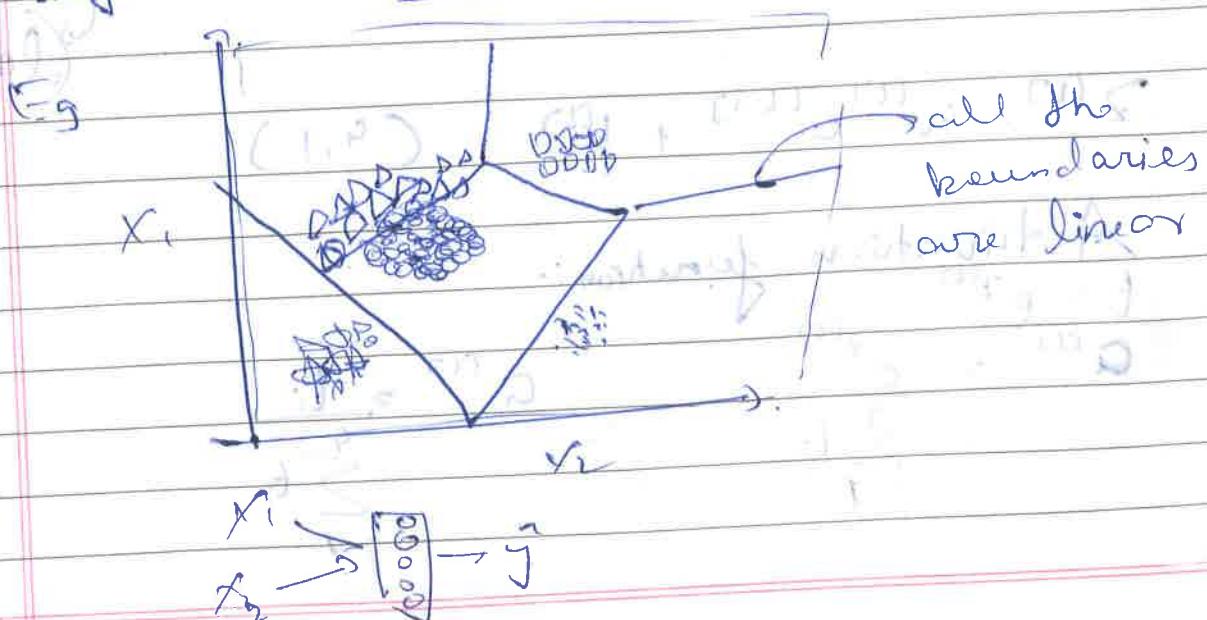
$$P(0|x) = e^{2(146.3)} = 0.042 = 4.2\%$$

$$P(1|x) = e^{-1/196.3} = 0.002 = 0.2\%$$

$$P(2|x) = e^{-2/196.3} = 0.998 = 11.4\%$$

$$P(3|X) = e^3 / 196.3 \approx 0.114 = 11.4\%$$

$\Rightarrow$  Decision boundary that we get through softmax is linear.



→ Training Softmax classifier

→ Softmax regression is the generalization of logistic regression to  $C$  classes.

→ Hardmax is just telling  $(1/0)$  & Softmax is telling the probabilities.

Eg.

$$\begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Softmax  $\rightarrow$  Hardmax

This means that if we reduce Softmax regression to just two classes it connects to Logistic regression.

Eg:  $\begin{bmatrix} 0.9 \\ 0.1 \end{bmatrix}$  as both the probabilities

add up to one so we don't need to calculate both of them, even calculating one of them would suffice. So while doing so it is reduced to Logistic regression. Don't go into the detailed math of it. It just reduces to that.

$$\{p, 1-p\}$$

$\Rightarrow$  Loss function:

$$L(\hat{y}, y) = -\sum_{j=1}^n (y_j) \log(\hat{y}_j)$$

so in this example:

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \hat{y} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$\begin{aligned} L(\hat{y}, y) &= -[(0) \log(0.3) + (1) \log(0.2) + (0) \log(0.1) \\ &\quad + (0) \log(0.4)] \\ &= -\log(0.2) \end{aligned}$$

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

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

now as we want  $L(\hat{y}, y)$  to be small as  $\hat{y}$  lies from 0 to 1 so we would want  $\hat{y}$  to be as close to 1 as possible to minimise loss function.

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

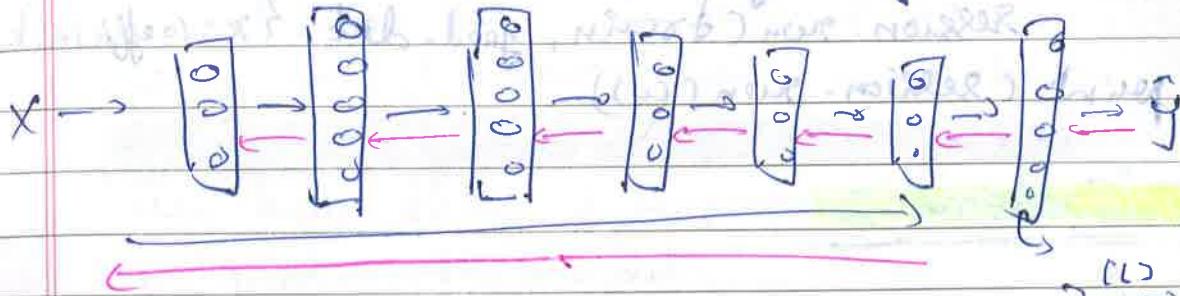
$$= \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \end{bmatrix}_{(4, m)}$$

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

$$= \begin{bmatrix} 0.2 \\ 0.3 \\ 0.1 \\ 0.4 \end{bmatrix}$$

(4, m)

⇒ Gradient descent with softmax :-



$$z^{(1)} \rightarrow a^{(1)} = \hat{y}^{(1)}$$

$$(4, 1)$$

$$z^{(2)} \rightarrow a^{(2)} = \hat{y}^{(2)}$$

$$(4, 1)$$

Backpropagation:-

$$\frac{dZ^{(1)}}{dZ^{(2)}} = \hat{y} - y$$

$$\frac{dZ^{(2)}}{dZ^{(1)}} = \hat{y} - y$$

→ Tensor Flow

Code example :-

```
import numpy as np
```

```
import tensorflow as tf
```

Coefficients = np.array([[1], [-20], [25]])

```
w = tf.Variable([0], dtype = tf.float32)
```

```
x = tf.placeholder(tf.float32, [3, 1])
```

```
cost = x[0][0]*w**2 + x[1][0]*w + x[2][0]
```

# (w - s) \*\* 2

```
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
```

```

init = tf.global_variables_initializer()
session = tf.Session()
session.run(init)
print(session.run(w))

```

with `tf.Session()` as `sess`

`session.run(w)`  
`print(sess.run(w))`

```

for i in range(1000):
    session.run(tf.assign, feed_dict={x: coefficients})
    print(session.run(w))

```

## Specialization 3

week 1:

### → Orthogonalization

This means that we should know that which hyperparameters would effect which aspect of the model. So that we can make a good model.

- ⇒ Chain of assumptions
- Fit training set well on cost function
  - if we (bigger) network,
  - if we use Adam optimization

$([(25, 25, 5), (100, 100, 10)])$  more q.v. - diff.

Fit dev set well on cost function

use  $\alpha$  =  $\alpha$  (regularization,  $\alpha$ )

$([(100, 100, 10)])$  Bigger dev set

$([(100, 100, 10)]) + \alpha \cdot \text{reg}(\text{cost})$

Fit test set well on cost function

$([(100, 100, 10)])$  (Bigger dev set)

↓  
Performs well in real world

PAGE NO. / / /  
DATE / / /

(change dev set or cost function)

They do not use "Early stopping" because this would effect both fitting on training set and dev set, so this is against orthogonalization as it states that it is better to use some those hyperparameters that tune only one part of the model.

- Single number evaluation metric are classified as cat and
- ⇒ Precision :- ( $\gamma$  of examples that are actually cats over the examples which are classified as cats)
- ⇒ Recall :- ( $\gamma$  of examples that are actually cat over the examples which are actually cat)

⇒ F1 score :- Harmonic mean of precision and recall :- 
$$\left( \frac{2}{\frac{1}{P} + \frac{1}{R}} \right) \Rightarrow \left( \frac{2PR}{P+R} \right)$$

- Satisfying and optimizing metric

| Classifier | Accuracy | Running time |
|------------|----------|--------------|
| A          | 90%      | 80ms         |
| B          | 92%      | 95ms         |
| C          | 95%      | 1.500ms      |

Optimizing → satisfying

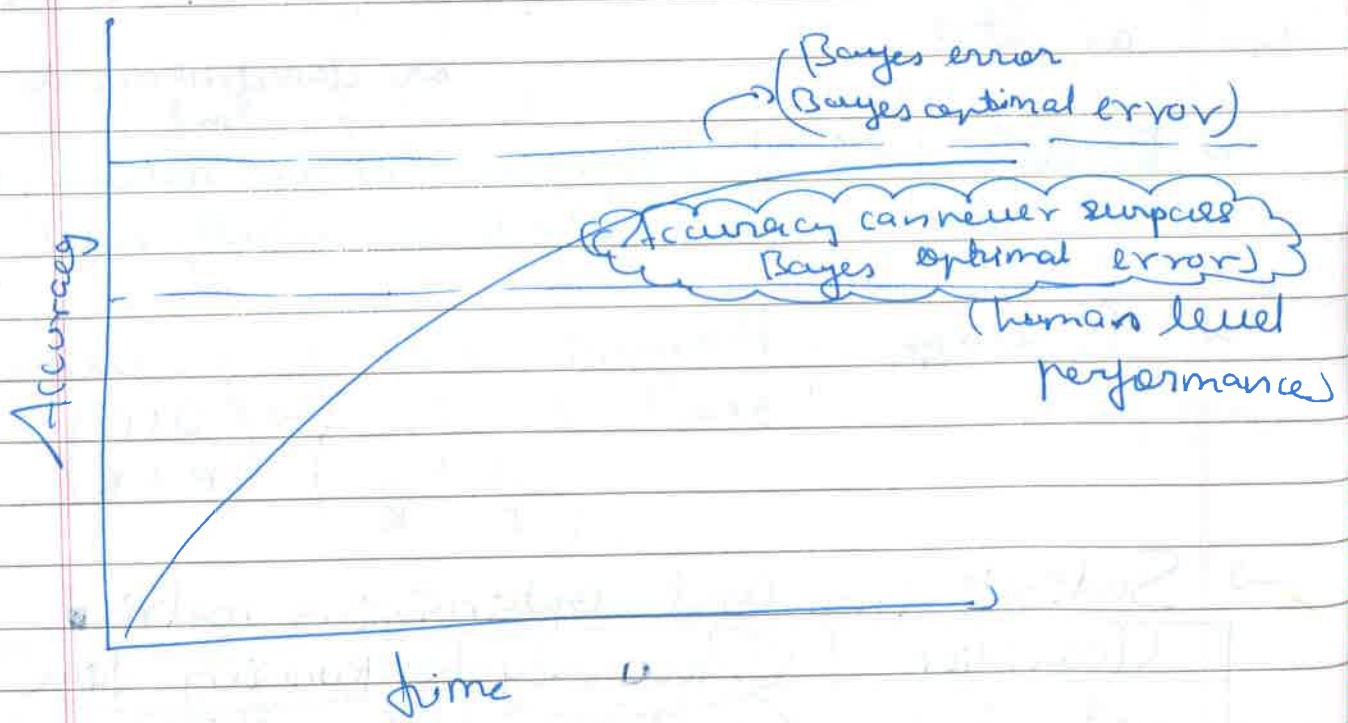
we could do:-

~~we could do following things~~

satisfy metrics:- That we could accept the model if it satisfies the threshold minimum condition.

We could in this case choose that the satisfy condition in Running time is that it should be less than 100ms.

- Train / dev / test distributions.  
(Dev set) & (test set) should come from the same distribution
- Comparing to human - level performance



If the accuracy of the model is below human level performance then there is a scope of improvement.

→ Avoidable bias.

Eg1 Recognizing cat images

Human error 1%

Training error 8%

Dev error 10%

(Focus on bias reduction)

Eg2 Recognizing cat images (Very low resolution)

Human error 7.5%

Training error 8%

Dev error 10%

(Focus on variance reduction)

→ avoidable bias :- 7%

Human error

→ avoidable bias :- 0.5%

→ Human level error

→ medical image classification example :-

(a) Typical human - 3%

(b) Typical doctor - 1%

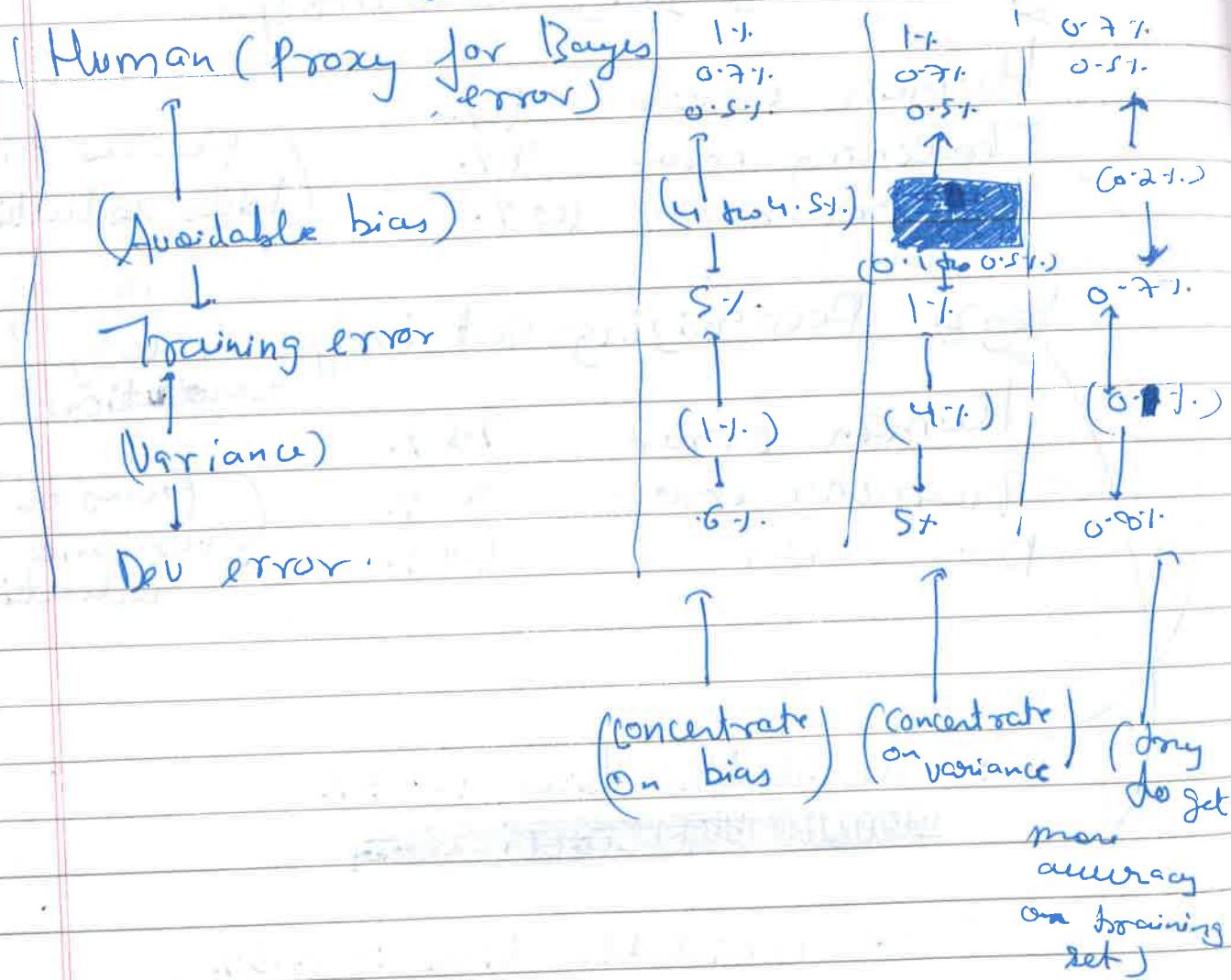
(c) Experienced doctor - 0.7%

(d) Team of experienced doctor - 0.5%

Now Bayes error  $\leq 0.5\%$ .

we might consider human level error to be  $0.5\%$  to check the model.

$\Rightarrow$  Error analysis example.



If we have an estimate on human level error then we can estimate on avoidable bias, variance & hence make decision that on which thing should we work on.

and use it proxy for bayes level error.

→ Surpassing human level performance

Eg)

Team of humans

- One human
- training error
- Dev errors

| Error (Eg 1) | Error (Eg 2) |
|--------------|--------------|
| 0.5%         | 0.8%         |
| 1%           | 1%           |
| 0.6%         | 0.3%         |
| 0.8%         | 0.4%         |

Eg where ML significantly surpasses human-level performance :-

- Online advertisement
- product recommendation
- logistics (predicting transit time)
- loan approvals.

in some cases ML has also surpassed human level performance:-

These cases are:-

- Medical diagnosis (Skin cancer etc)
- Speech recognition
- Computer vision



Week 9

## → Carrying out error analysis

→ Example that you are building a model that would recognize cat & tell whether that there is a cat in the picture

So now after building the initial model we could analyze the ~~the~~ errors

Ideas for cat detection:-

- 1) Fix pictures of dogs being recognized as cats
- 2) Fix great cats (lions, panthers, etc.) being misrecognized
- 3) Improve performance on blurry images.
- 4) misclassifying Instagram filters as bats

| Images | Dog | Great cats | Blurry | Instagram | Correct |
|--------|-----|------------|--------|-----------|---------|
| 1      | ✓   |            |        |           | Pitbull |
| 2      |     |            | ✓      | ✓         |         |
| ?      |     | ✓          | ✓      |           |         |
| ?      | ?   |            |        |           |         |
| Total  | 81  | 437        | 617    | 121       | 1       |

→ A table like this is made on error images so that we can create categories our scaling problems to make our model better.

→ Cleaning up incorrectly labelled data  
This is the dataset that is wrong, in which the data that we are inputting into the model for training are incorrectly

It is more important to correct incorrectly labelled data in dev/test set as compared to learning set.

## labelled

Sometimes it is fine to train ML models on mislabelled data when the error is random.

The problem arises when the input that feed in is incorrectly labelled systematically. Eg. when all white dogs are incorrectly labelled as cats

So we make a table on all  labelled data.

| Image    | Dog | Correct cat | Blurry | Incorrectly labelled | Comments |
|----------|-----|-------------|--------|----------------------|----------|
| 99       |     |             |        |                      |          |
| 99       |     |             |        | ✓                    |          |
| 100      |     |             | ✓      |                      |          |
| ...      | ... | ...         | ...    | ...                  |          |
| 1 of 100 | 8.1 | 43.1        | 61.1   | 6.1                  |          |

Here overall dev set error :- 10%.

Error due to incorrect labels - 0.6%.

Errors due to other causes - 9.4%.

So there is not much need to concentrate on incorrectly labelled data now.

But suppose you have an overall dev set error of 10%. Then there is need to correct the incorrectly labelled data.

⇒ Correcting incorrect dev / test set examples

- ) Apply same process to your dev and test set to make sure they continue to come from the same distribution.
- ) Consider examining examples your algorithm got right as well as ones it got wrong.



⇒ It is important that we make sure (dev) & (test) set are from the same distribution even though training set is not so.

→ A tip:-

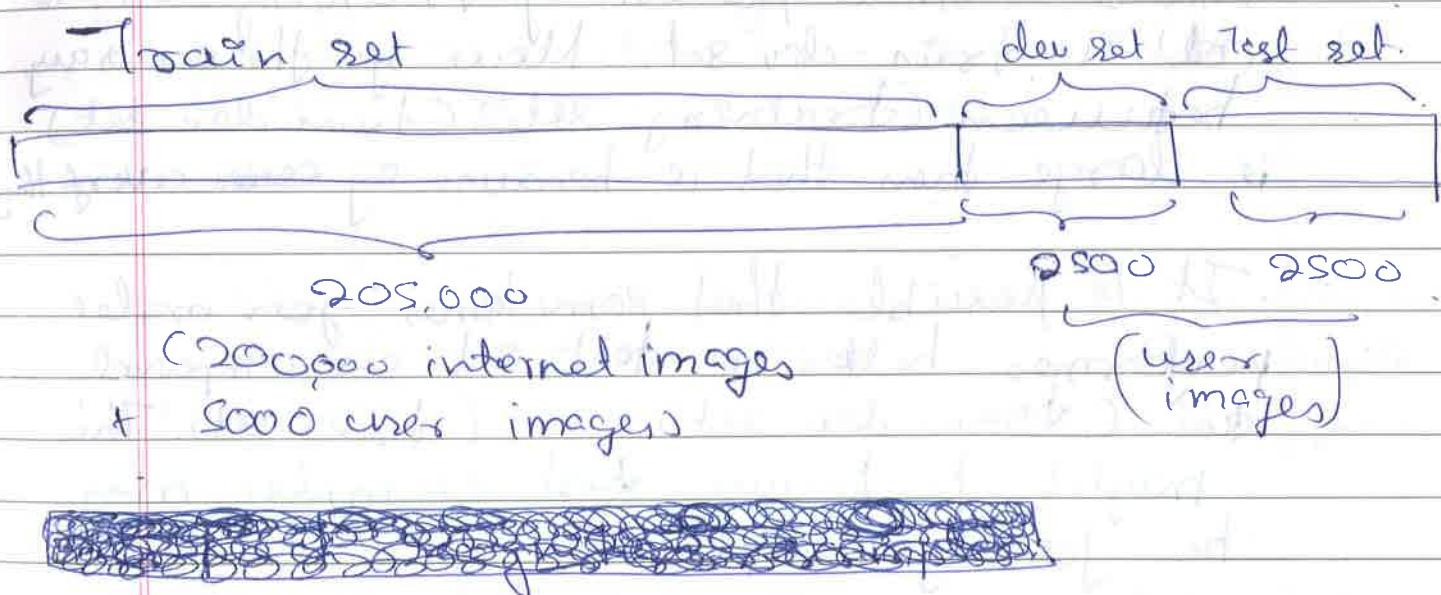
- ⇒ Set up dev / test set and metric
- ⇒ Build initial system quickly
- ⇒ Use Bias / Variance analysis and error analysis to prioritize next steps.

Guideline:- Build your first system quickly, then iterate

→ Training & testing on different distributions

- ⇒ Let's say that you have to make a cat recognizer model.  
So for testing user would input data (images of cat) clicked from their phone.  
But for training you have 200,000 images of  cat taken from the internet (professionally clicked with much

better quality). You do a survey and ask people to click cat photos and send it to you, like this you get 10,000 images. Now these images come from the same distribution as the test set but the dilemma is that you have 200,000 images of cat that don't come from the same distribution as the test set so the most preferable option should be the following:



→ Bias  $\propto$  Variance with mismatched distribution.

Human level

{ avoidable bias }

Training set error  
Variance

Train - dev set error  
} (data mismatch)

Dev error

{ degree of overfitting is coming from some distribution }

Test error

→ Sometimes as we have seen the distribution of dev set & testing set might be fairly different so if we train our model on training set and then test it on dev set, and if we get low accuracy on dev set then we may not know that the lower accuracy is due to overfitting or data mismatch.

→ To solve this problem we randomly remove a small portion of training set and call it train-dev set. Now if the accuracy between (training set) & (train-dev set) is large then that is because of ~~overfitting~~ overfitting.

⇒ It is possible that sometimes your model performs better on test set as compared to (train-dev set) and (train set). This might be because test examples may be fairly simple.

→ Generating artificial data.

⇒ Eg you need to generate artificial data of people in car for a gear view matrix speech recognition model. And you have 10,000 hours of normal talking data of people in quiet environment & 1 hour of car noise data.  
↳ (just)

What you could do is that shuffle that 1 hour of just car noise data and overlap with 10,000 hours of quiet data to generate artificial data.

This would pose a problem though.

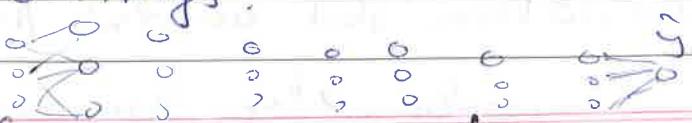
Our model may convert to that 1 hour of car sounds. For normal human ears, all of 10,000 hours of car noise is unique.

→ and 1 hour of car noise overlapped 10,000 times might sound the same but the computer may convert to just that 1 hour of car noise.

### → Transfer learning

Transfer learning is when we use one machine learning model to solve another problem. For example you built a model for image detection. Now you want to build a model for radiology X-ray detection. Now the basis of radiology & image recognition is same, like boundary recognition, shape recognition. So if we'll make sure to train radiology model on image recognition model after adding some additional neural layers on the existing neural network.

And also this would make sure when we have just 100-1000 images for radiology detection but 10,000,000 examples for image detection because we can't train a radiology model solely on the basis of 100 images.



→ (Shared layers for image detection) (Additional layers for radiology detection)

When does transfer learning make sense:

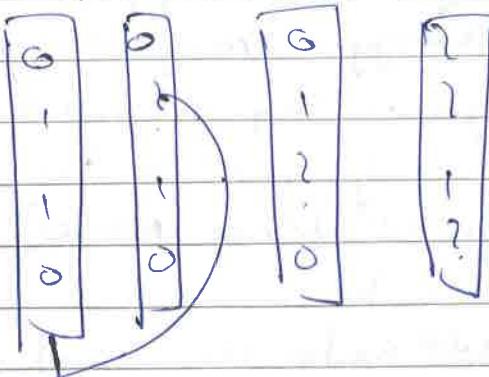
- (i) Task A & Task B have the same inputs  $\times$
- (ii) You have more data on Task A as compared to Task B
- (iii) Low level features from A could be helpful for learning B.

→ Multi-task learning.

Suppose we have to make a model that would give us an output like:

- $G \rightarrow P(1)$  for traffic light
- $O \rightarrow P(1)$  for stop sign
- $O \rightarrow P(1)$  for humans on road
- $O \rightarrow P(1)$  for trucks on road

so this would mean:



(no traffic light, there is stop signal,  $\times$ )

there are humans and there is no truck

→ (This would mean that we don't know whether there is a traffic stop sign or not.)

$$\text{Loss: } \hat{y}^{(i)} \rightarrow \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^4 L(\hat{y}_j^{(i)}, y_j^{(i)})$$

$$- (y_j^{(i)} \log \hat{y}_j^{(i)} - (1-y_j^{(i)}) \log (1-\hat{y}_j^{(i)}))$$

So the other way to solve this problem could have been to make 4 different neural network models for 4 of these tasks. But it is better to make a single neural network for multiple sub-tasks because the initial layers for all the four tasks would be same and we could combine all the dataset for these 4 tasks for a single dataset for one neural network. So one combined neural network would have a much larger dataset than the dataset available.

Currently multi-tasking is being used in computer vision.

Still today transfer learning is used more than multi-task learning.

- End-to-End deep learning.
  - Initially we end to break the learning process into different parts. For example if for a voice recognition system we would first make a model to recognize phonemes (sound like "a", "p", "f", etc) then further on.

Because as this the algorithm many times does not have a free hand to go about this restricting its performance. But for this end to end deep learning we need large amounts of data. ~~data~~. If we have less data then it would be better to manually guide the model at steps.

## Specialization 4 week 1

### Convolution neural networks

→ Vertical Edge detection

The square would give the value

(Convolution)

The square would give this value

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| 3 | 0 | 1 | 2 | 7 | 4 |
| 5 | S | 0 | 9 | 3 | 1 |
| 2 | 7 | 2 | 5 | 1 | 2 |
| 0 | 1 | 3 | 1 | 7 | 8 |
| 4 | 2 | 1 | 6 | 2 | 8 |
| 2 | 4 | 5 | 2 | 3 | 9 |

|    |   |    |
|----|---|----|
| -1 | 0 | 1  |
| 1  | 0 | -1 |
| -1 | 0 | 1  |

4 x 4

(filter)

(gray scale image)

Looking at the first  $(3 \times 3)$  square

|   |   |   |    |
|---|---|---|----|
| 1 | 3 | 0 | -1 |
| 1 | 1 | 5 | 7  |
| 1 | 2 | 0 | 2  |

$$\begin{aligned}
 & \Rightarrow (3)(1) + (1)(1) + (2)(1) + (0)(0) \\
 & + (0)(5) + (0)(7) + (-1)(1) + (-1)(2) \\
 & + (-1)(2) \\
 \Rightarrow & 3 + 1 + 2 - 1 - 8 - 2 \\
 = & -5
 \end{aligned}$$

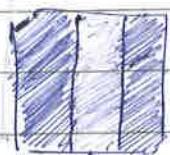
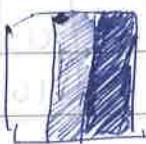
Ex1

|    |    |    |   |   |   |
|----|----|----|---|---|---|
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |

 $6 \times 6$ 

\*

|   |   |   |   |    |    |   |
|---|---|---|---|----|----|---|
| 1 | 0 | 1 | 0 | 30 | 30 | 0 |
| 1 | 0 | 1 | 0 | 30 | 30 | 0 |
| 1 | 0 | 1 | 0 | 30 | 30 | 0 |

 $8 \times 6$  $4 \times 4$ 

→ edge in the picture

we have  
detected  
the output edge.

Fig 2: Now if we change our initial matrix  
to:-

|   |   |   |     |     |   |
|---|---|---|-----|-----|---|
| 0 | 0 | 0 | 70  | 0   | 0 |
| 0 | 0 | 0 | 0   | 100 | 0 |
| 0 | 0 | 0 | 100 | 100 | 0 |
| 0 | 0 | 0 | 100 | 0   | 0 |
| 0 | 0 | 0 | 100 | 0   | 0 |

 $(6 \times 6)$ 

\* Convolution  
(3x3)

|   |      |     |   |
|---|------|-----|---|
| 0 | 1-30 | -30 | 0 |
| 0 | -30  | -30 | 0 |
| 0 | -30  | -30 | 0 |
| 0 | -30  | -30 | 0 |

→ we will  
have -30 instead  
of +30.

→ horizontal edge detector:-

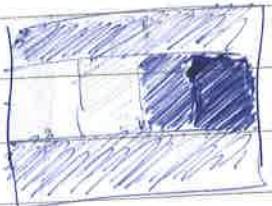
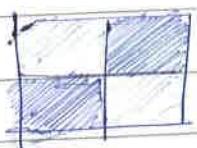
|    |    |    |    |
|----|----|----|----|
| 1  | 1  | 1  | 1  |
| 0  | 0  | 0  | 0  |
| -1 | -1 | -1 | -1 |

Ex

|    |    |    |    |    |    |
|----|----|----|----|----|----|
| 10 | 10 | 10 | 0  | 0  | 0  |
| 10 | 10 | 10 | 0  | 0  | 0  |
| 10 | 10 | 10 | 0  | 0  | 0  |
| 0  | 0  | 0  | 10 | 10 | 10 |
| 0  | 0  | 0  | 10 | 10 | 10 |
| 0  | 0  | 0  | 10 | 10 | 10 |

|   |    |    |    |
|---|----|----|----|
| * | 1  | 1  | 1  |
|   | 0  | 0  | 0  |
|   | -1 | -1 | -1 |

|    |    |     |     |
|----|----|-----|-----|
| 0  | 0  | 0   | 0   |
| 30 | 10 | -10 | -30 |
| 30 | 10 | -10 | -30 |
| 0  | 0  | 0   | 0   |



⇒ Now we could choose different filters instead of  $\begin{bmatrix} 1 & 0 & -1 \\ 1 & 3 & -1 \\ 1 & 2 & -1 \end{bmatrix}$ .

Some of the other popular filters are:

|   |   |    |
|---|---|----|
| 1 | 0 | -1 |
| 2 | 0 | -2 |
| 1 | 0 | -1 |

Sobel filter

|    |   |     |
|----|---|-----|
| 3  | 0 | -3  |
| 10 | 0 | -10 |
| -3 | 0 | -3  |

Sobel filter



→ vertical

|    |    |    |
|----|----|----|
| 1  | 2  | 1  |
| 0  | 0  | 0  |
| -1 | -2 | -1 |

|    |     |    |
|----|-----|----|
| 3  | 0   | 3  |
| 0  | 0   | 0  |
| -3 | -10 | -3 |



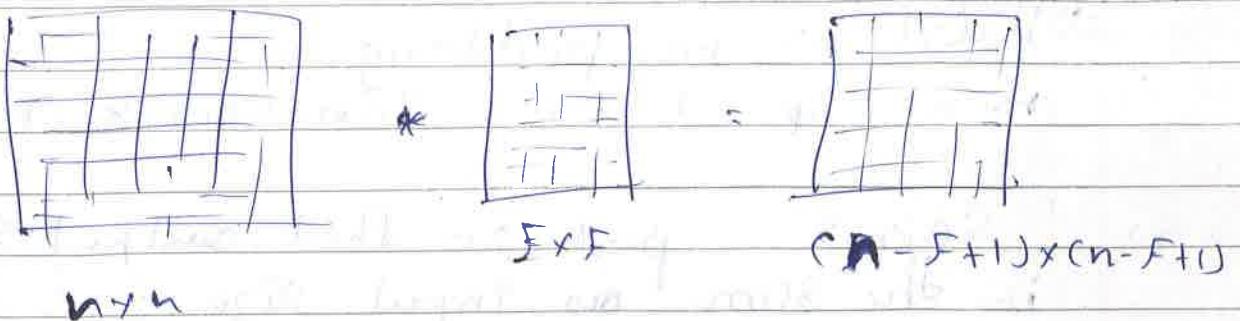
→ horizontal

⇒ We could also learn these as parameters:

|       |       |       |
|-------|-------|-------|
| $w_1$ | $w_2$ | $w_3$ |
| $w_4$ | $w_5$ | $w_6$ |
| $w_7$ | $w_8$ | $w_9$ |

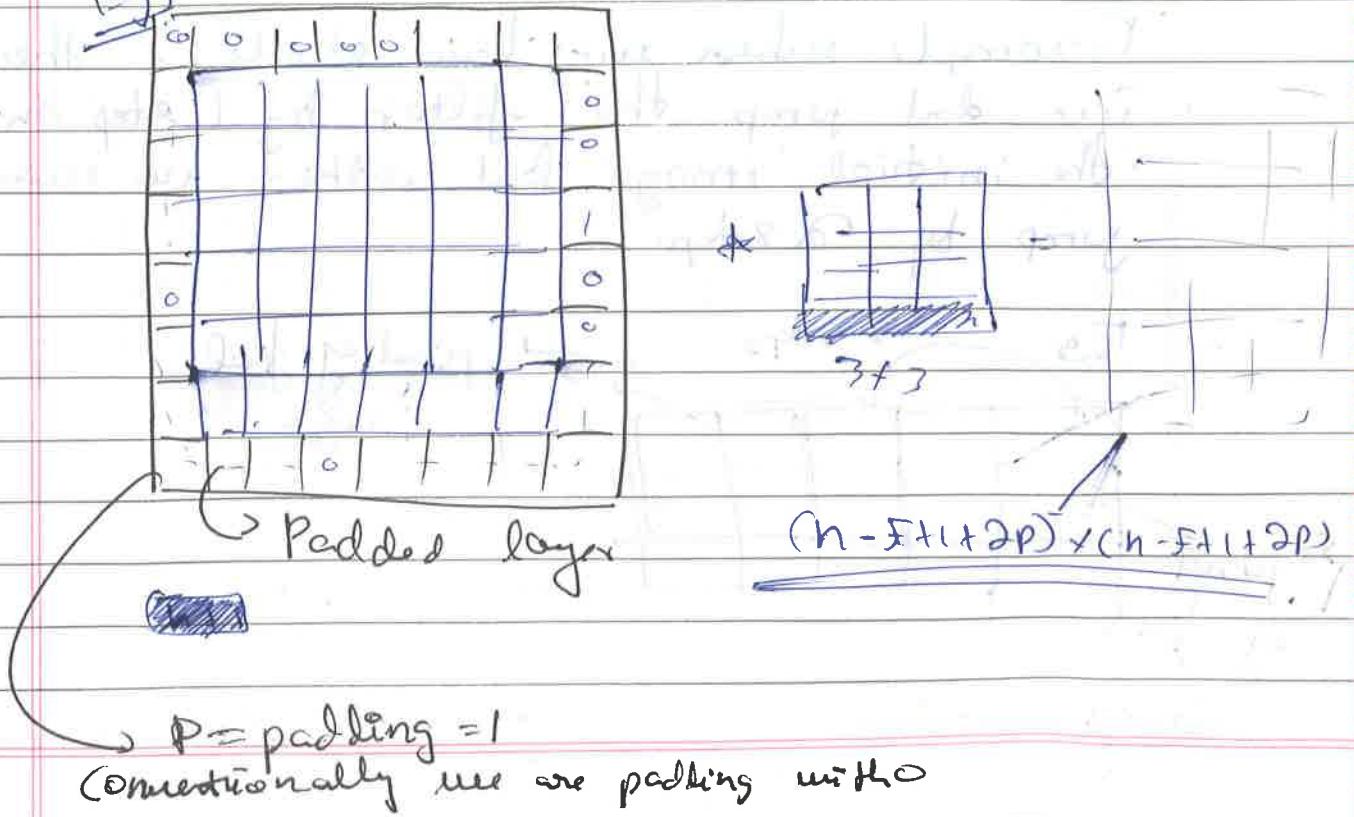
do not only detect vertical/horizontal edges but also detect edges at a certain angle.  
E.g.  $75^\circ$  edges.

## → Padding



- The disadvantage of this is that the size of the image (pixels) is reduced from  $(n \times n)$  to  $(n-F+1) \times (n-F+1)$
- One more disadvantage is that the pixels at the edge give relatively less info than pixels in the middle. This is because pixels in the middle are used more to determine the output image.
- To solve this we introduce padding. Padding is adding layers on the border of initial image.

Eg)



⇒ Valid and same convolutions.

⇒ "Valid" → no padding.

$$(n \times n) * (F \times F) = (n-F+1) \times (n-F+1)$$

⇒ "Same" : pad so that output size is the same as input size.

$$(n \times n) * (F \times F) = (n-F+1+2p) * (n-F+1+2p)$$

we need:

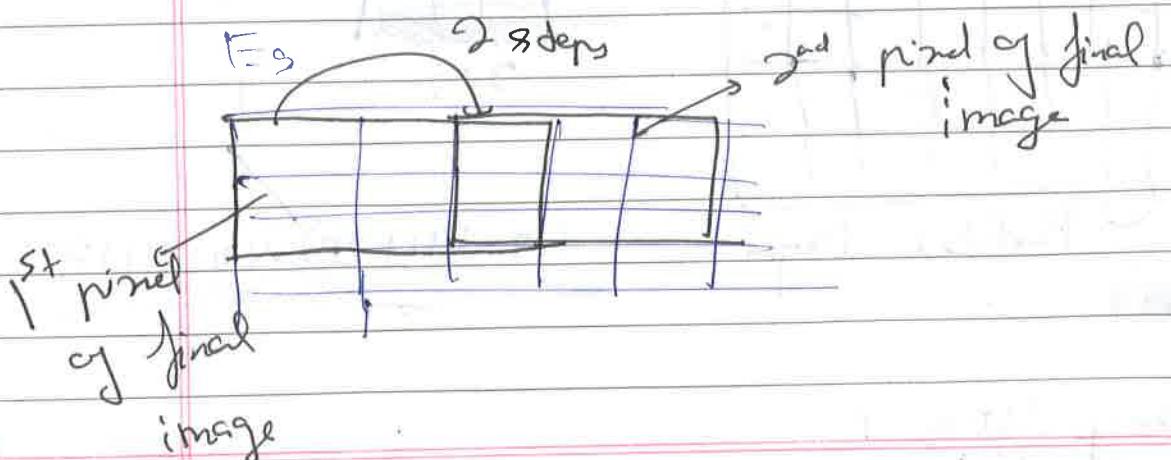
$$n = n-F+1+2p$$

$$\Rightarrow 2p = F-1 \Rightarrow p = \frac{F-1}{2}$$

F is almost always odd. We rarely see F as even.

⇒ Strided convolution.

For example when we have stride = 2 then we don't jump the filter by 1 step on the initial image but rather we should jump by 2 steps.



Let  $s$  be the number of strides. Then:

$$\rightarrow (n \times n) * (F \times F) = \left[ \frac{n+2p-F+1}{s} \right] \times \left[ \frac{n+2p-F+1}{s} \right]$$

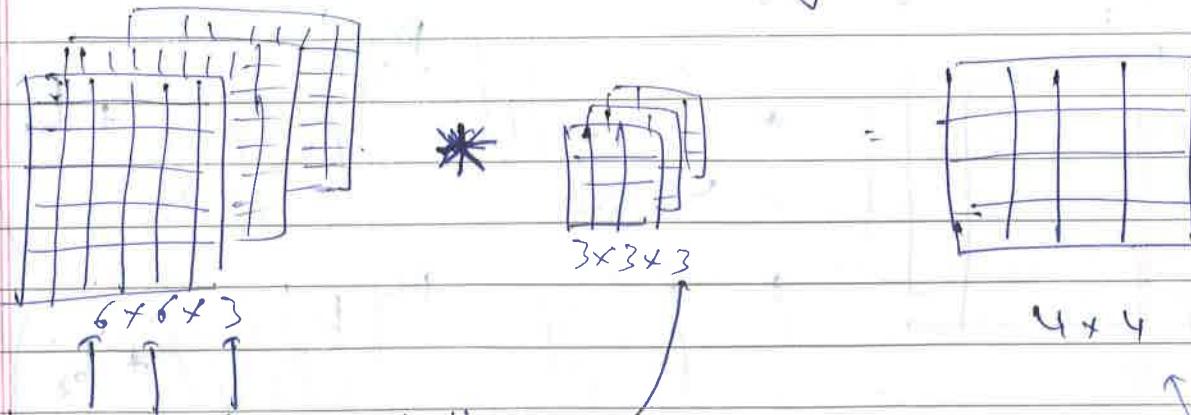
$\hookrightarrow$  G.I.F

Note:- What we did, what we call convolution is called cross-correlation in mathematics. But we call it convolution in deep learning. Convolution in mathematics is little different. It is basically rotating the filter  $180^\circ$  clockwise and then applying it. 

(or you could just flip it vertically then horizontally)

$\rightarrow$  Convolutions on volumes.

$\rightarrow$  Convolutions on RGB image.



height width (number of channels)

(Other dimensions must be equal)

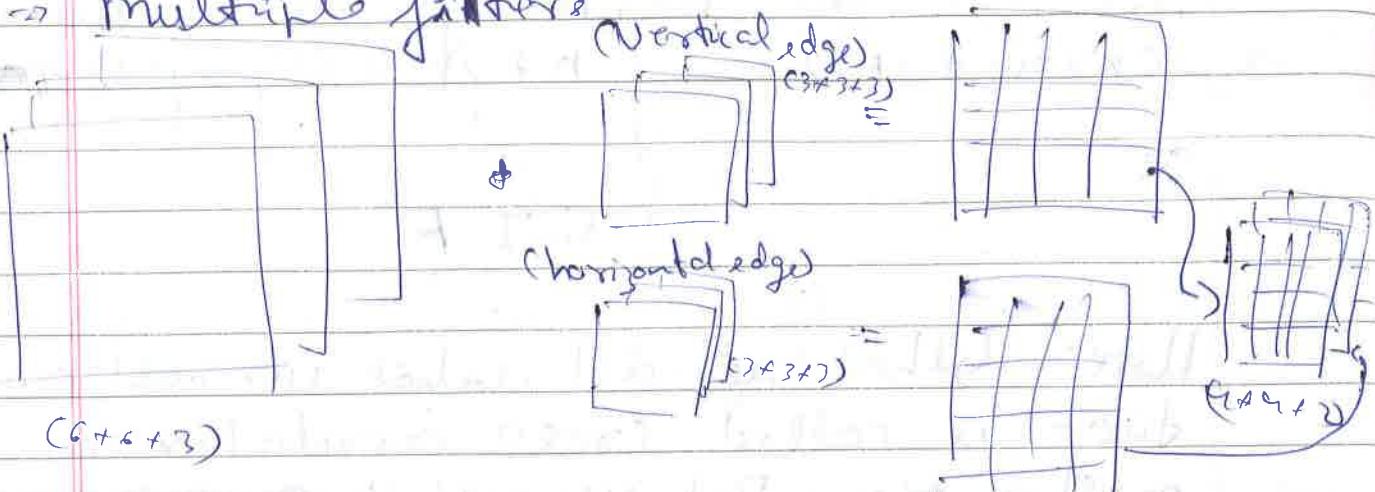
$\hookrightarrow$  This would lead to:

$(p=1, 0)$



$(4 \times 4 \times 3)$  then we need to add the weight along 2 direction (+) & bias term

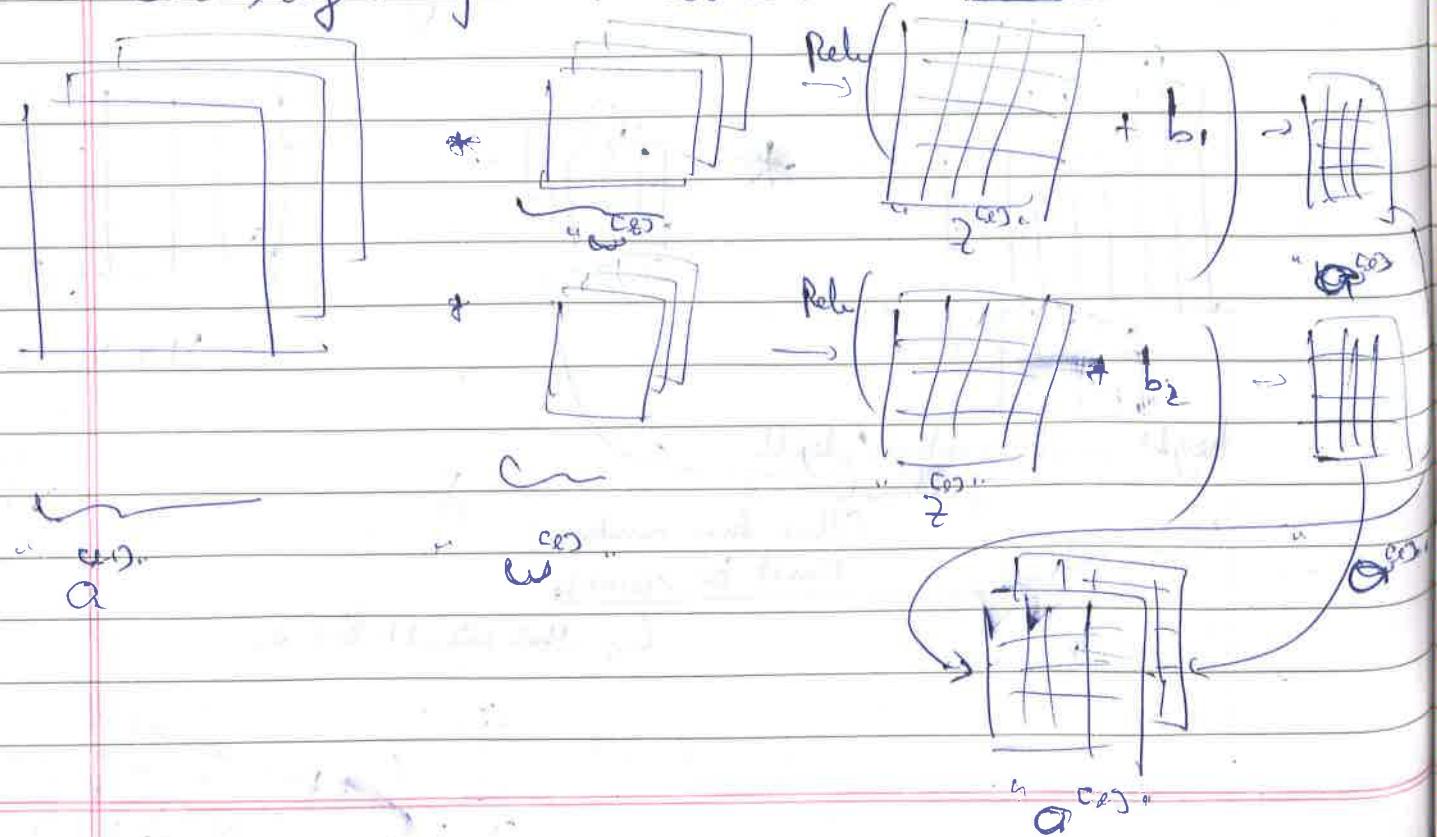
→ multiple filters



Summary

$$(n \times n \times n_c)^{\text{number of channels}} \times (F \times F \times n_c) \rightarrow (n - F + 1) \times (n - F + 1) \times n_c^{\text{number of filters}}$$

→ One layer of a convolutional ~~network~~ network.



~~So what is the total number of parameters in one layer are:-~~

So the total number of parameters when we have 16 filters are:-

$$16 \times (3 \cdot (F \cdot F) + 1) \Rightarrow 30F^2 + 10$$

$$\text{let } F = 3$$

$$\text{thus } \Rightarrow 30(9) + 10 = \underline{\underline{280}}$$

So now the total number of parameters in one layer is 280 even if we have an image that has a lot of pixels.

⇒ Summary of notation:-

If  $l$  is a convolution layer:-

$F^{(l)}$  = filter size

$P^{(l)}$  = padding

$S^{(l)}$  = stride

$n^{(l)}$  = number of filters

(each filter is  $= F^{(l)} \times F^{(l)} \times n_c^{(l)}$ )

(Activations:  $A^{(l-1)} \rightarrow n_h^{(l)} \times n_w^{(l)} \times n_c^{(l)}$ )

$(A^{(l-1)} \rightarrow M \times n_h^{(l)} \times n_w^{(l)} \times n_c^{(l)})$

(Weights:  $n_h^{(l)} \times F^{(l)} \times F^{(l)} \times n_c^{(l)}$ )

number of filters.

(Input:  $n_h^{(l-1)} \times n_w^{(l-1)} \times n_c^{(l-1)}$ )

(output:  $n_h^{(l)} \times n_w^{(l)} \times n_c^{(l)}$ )

$n_h^{(l)} = \left\lceil \frac{n_h^{(l-1)} + 2P^{(l)} - F^{(l)} + 1}{S^{(l)}} \right\rceil$

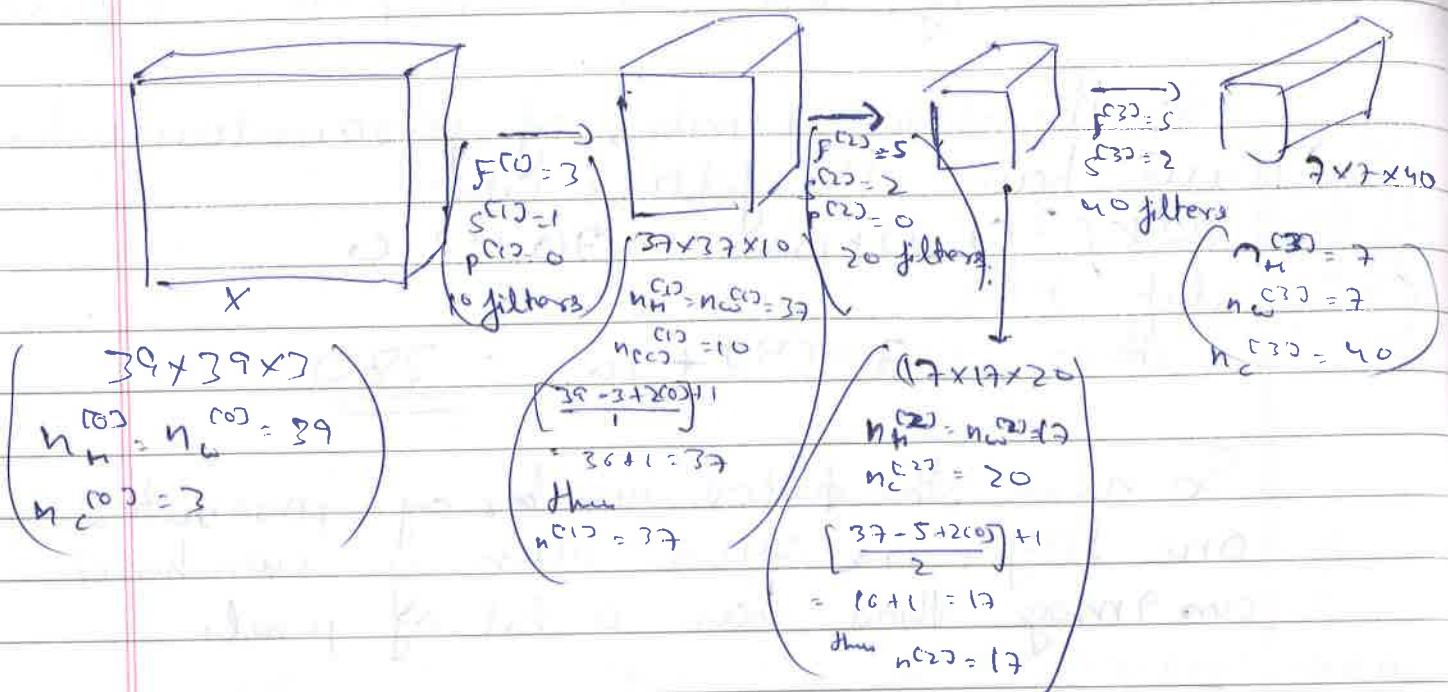
$n_w^{(l)} = \left\lceil \frac{n_w^{(l-1)} + 2P^{(l)} - F^{(l)} + 1}{S^{(l)}} \right\rceil$

no. of channels in the previous layer

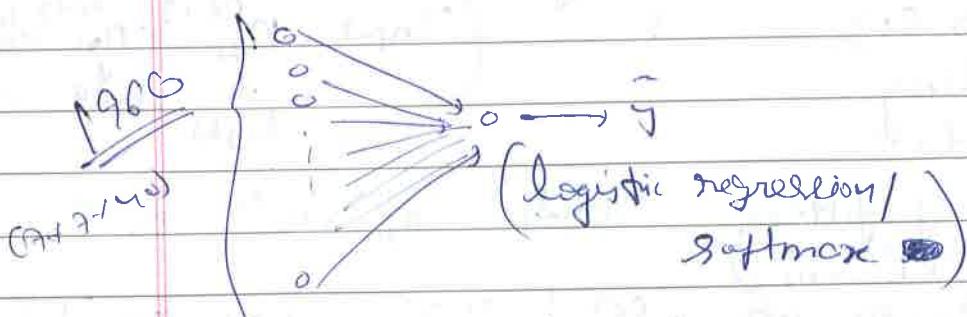
bias:  $n_c^{(l)}$  it would be more convenient if we write it as:  $(1, 1, 1, n_c^{(l)})$

$b$ , the constant which is added.

⇒ Example of convolution network

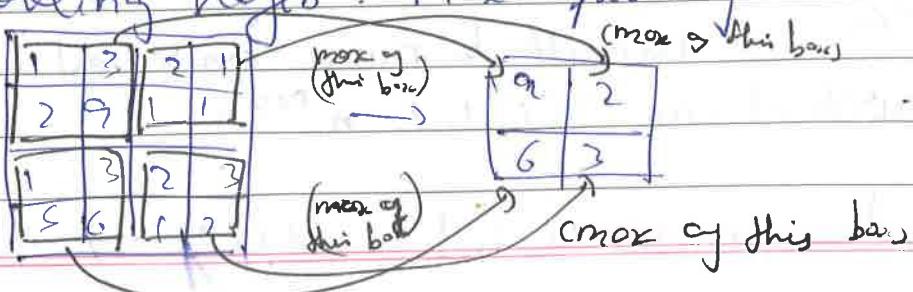


now total units we have are:  $7 \times 7 \times 40 = 1960$



- ⇒ Types of layers in a convolutional network
- Convolution (CONV)
  - Pooling (POOL)
  - Fully connected (FC)

⇒ Pooling layer: max pooling



You can use pooling layer to reduce  
now, now but not  $n_c$

PAGE NO.

DATE

hyper parameters:-  $F = 2$  (boxes are  $2 \times 2$ )  
 $s = 2$  (stride is 2)

Q) There are no parameters to learn in max pooling.  
Thus no gradient descent would be applied on any parameters of max pooling as there are none.

The intuition behind max pooling is that we would be able to detect in which quadrant or in which part of the image we have ~~any~~ a particular feature. For example an eye would have very high pixels. So in the ~~following~~ <sup>left</sup> example we had  $\rightarrow$  in the upper ~~right~~ quadrant, this we may ~~infer~~ <sup>infer</sup> that we have a cat's eye in the upper left quadrant.

→ Average pooling

Here instead of taking the max we take the average of all the numbers in that box.

|   |   |   |   |
|---|---|---|---|
| 1 | 3 | 2 | 1 |
| 2 | 9 | 1 | 1 |
| 1 | 4 | 2 | 3 |



|      |      |
|------|------|
| 3.75 | 1.25 |
| 4    | 2    |

When there is no padding:-

$$n_H = \left[ \frac{(n_H \text{ pixels} - F)}{\text{stride}} \right] + 1$$

$$n_W = \left[ \frac{(n_W \text{ pixels} - F)}{\text{stride}} \right] + 1$$

$$n_C = n_C \text{ prev.}$$

When there is padding  
add + 2p as usual

→ Hyperparameters:-

$F$ : filter size

$s$ : stride

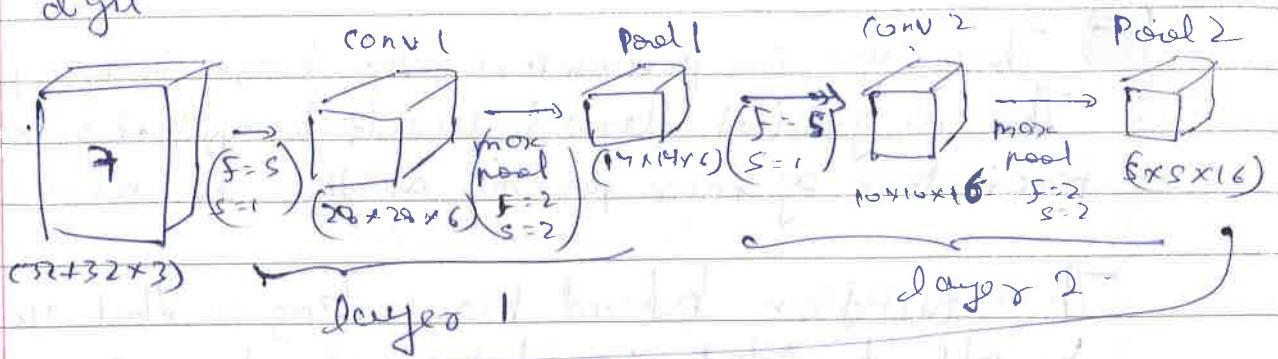
Max or average pooling

Usually we don't do padding here

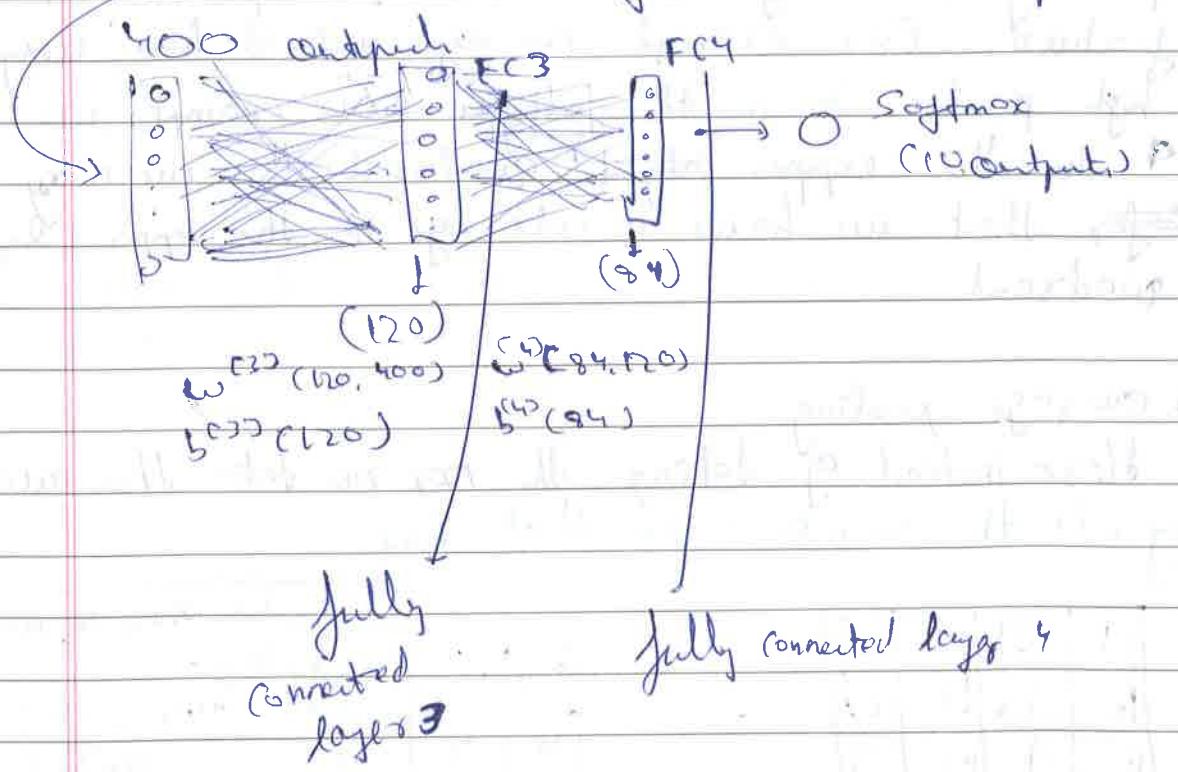
Input:  $n_H \times n_W \times n_C$

Output:  $\left[ \frac{n_H - F}{s} + 1 \right] \times \left[ \frac{n_W - F}{s} + 1 \right] \times n_C$

max pooling also performs as a noise suppressant. It discards the noisy activations altogether and performs de-noising along with dimensionality reduction on the other hand average pooling simply performs dimensionality reduction as a noise suppression mechanism. Hence, the CNN example, ~~max pooling~~ that more pooling performs let's build a neural network to identify the digit



Now we would flatten our output to



As we move forward notice that width, height would decrease & number of channels increases.

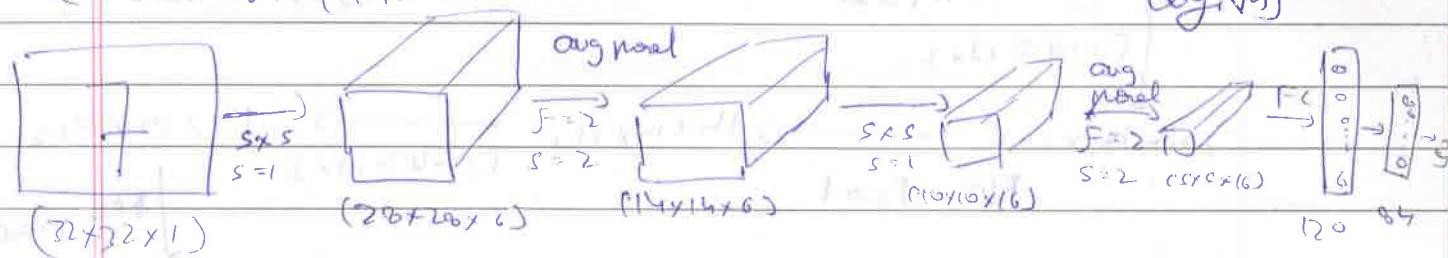
→ Convolution helps in parameter sharing.

A feature detector (such as a vertical edge detector) that is useful in one part of the image is probably useful in another part of the image.

- less better than average pooling
- sparsity of connections: In each layer, each output value depends only on a small number of inputs
- convolution neural network are less prone to regularization.
- (NN are immune to translation invariance. This means even if some pixels are shifted then also it would give good results.)

## Week 2 Specialization 4 (Case Studies)

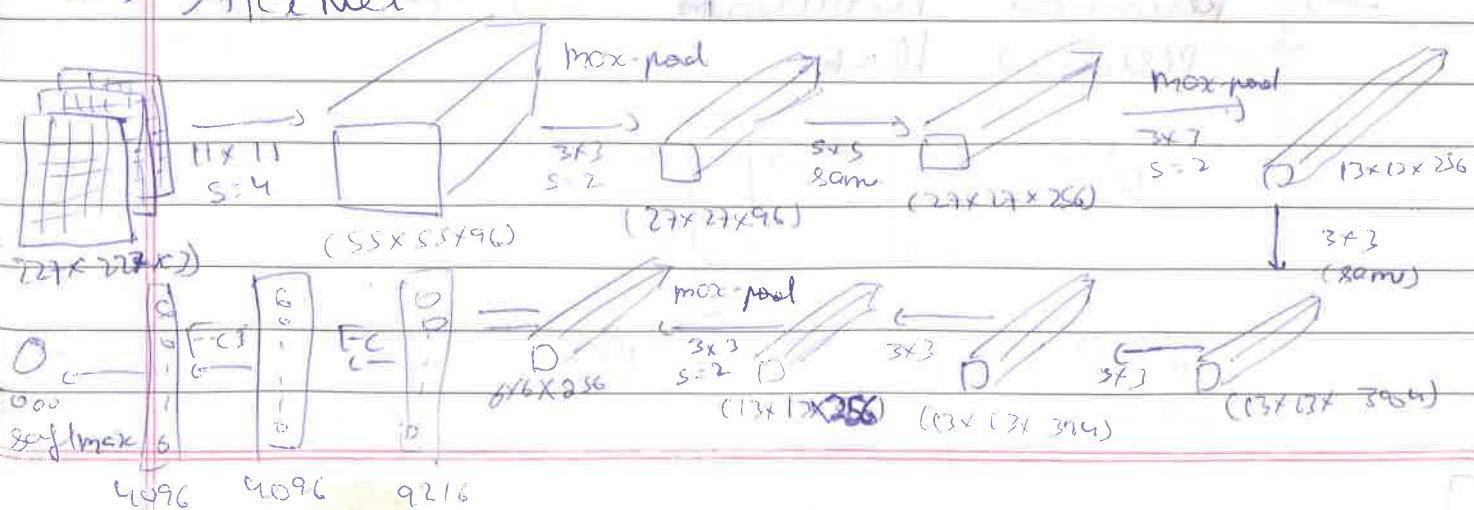
→ LeNet-5 (used to recognize handwritten digits)  
(made in 1990's)



This has around 60k parameters, but networks used today have millions of parameters.

Now as we move forward nn, nn & nt

→ AlexNet

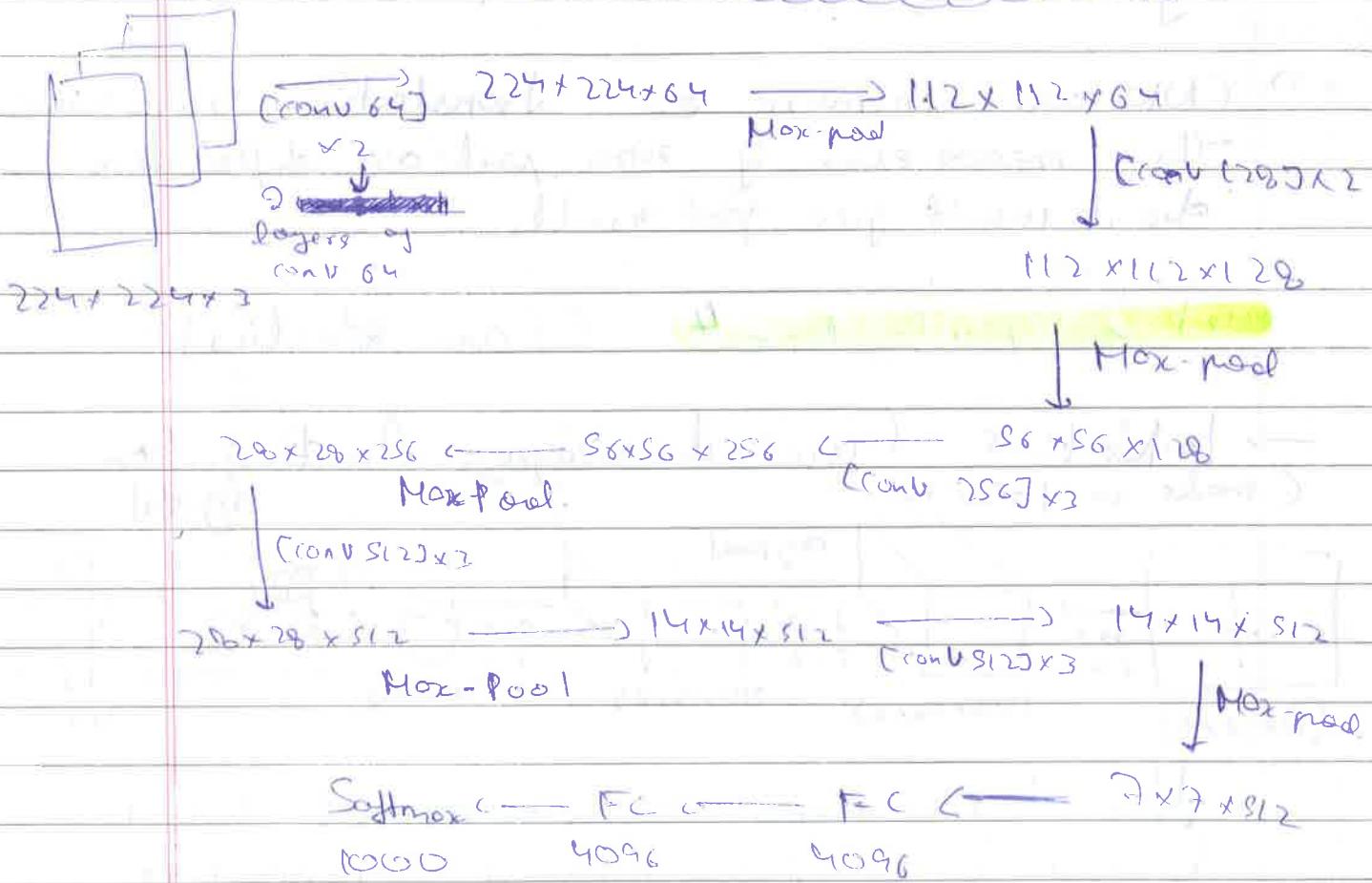


This is much bigger than LeNet-5, it has around 60 million parameters.

ReLU activation function was used here.

(Here Conv means  $3 \times 3$  filter,  $S=1$ , padding=1)

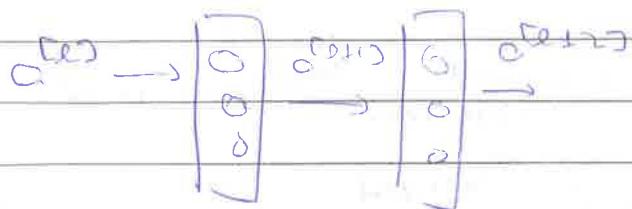
→ VGG-16 (and Max-Pool =  $2 \times 2$ ,  $S=2$ )



This has around 138 million parameters.

→ Residual Networks (ResNets)

→ Residual block



$$z^{0+12} = w^{0+12} z^{0+12} + b^{0+12}$$

$$z^{0+13} = g(z^{0+12})$$

Resnets prevent exploding or vanishing gradients

|          |     |
|----------|-----|
| PAGE NO. | 111 |
| DATE     |     |

$$z^{(l+2)} = w^{(l+2)} a^{(l+1)} + b^{(l+2)}$$

$$a^{(l+2)} = g(z^{(l+2)})$$

$$c^{(l+2)} \rightarrow \text{linear} \rightarrow \text{Relu} \rightarrow \text{linear} \rightarrow \text{Relu} \rightarrow c^{(l+2)}$$

This is the main path.

Note in residual block we will do:-

"short cut" / skip connection

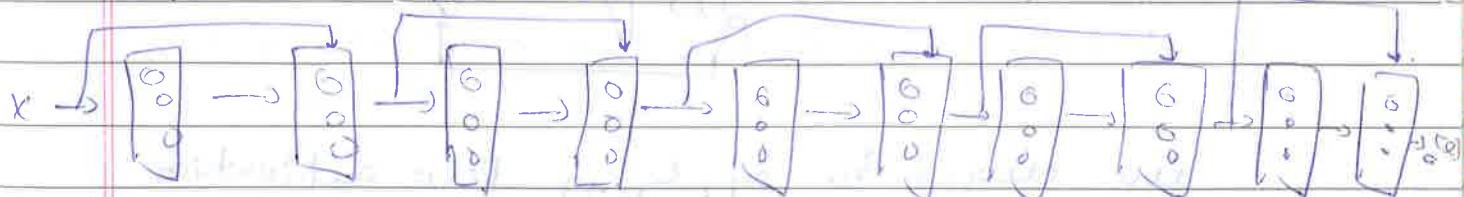
$$c^{(l+2)} \rightarrow \text{linear} \rightarrow \text{Relu} \xrightarrow{c^{(l+3)}} \text{linear} \rightarrow \text{Relu} \rightarrow c^{(l+2)}$$

"main path"

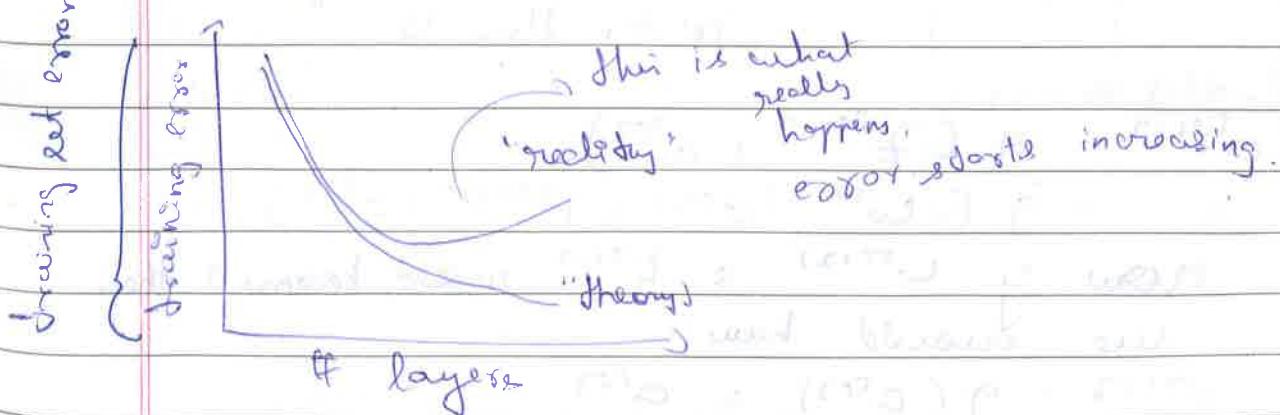
$$\text{then } c^{(l+2)} = g(c^{(l+2)} + c^{(l+3)})$$

Using residual blocks helps us to make  $\rightarrow$  much deeper neural ~~networks~~ networks  $\times$

$\Rightarrow$  Residual network

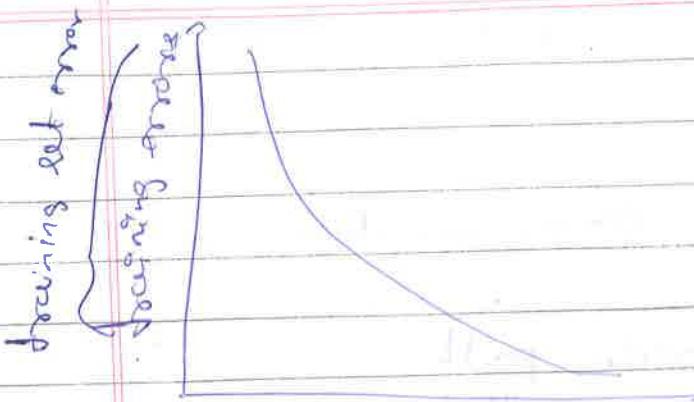


In plain neural network



→ In residual networks.

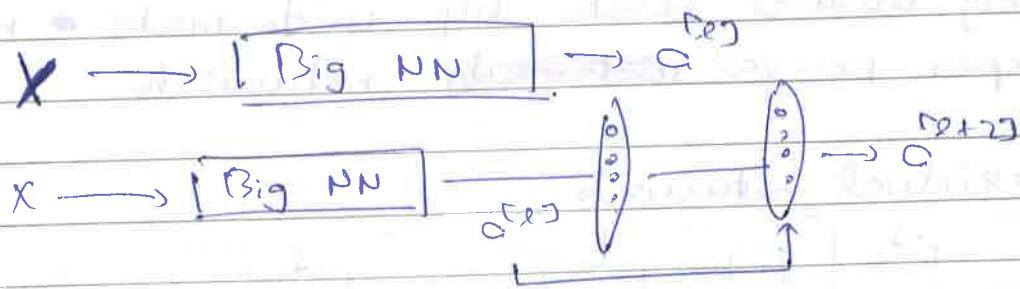
|          |     |
|----------|-----|
| PAGE No. | 111 |
| DATE     |     |



H layers.

Thus if we use residual networks then on increasing the ~~size~~ number of layers, the error keeps decreasing.

→ Let's ResNet work



we could be applying ReLU activation function thus  $a \geq 0$

$\rightarrow c$  would always be greater than 0.

$$\begin{aligned} c^{L_{l+1}} &= g(z^{L_{l+1}} + a^{L_l}) \\ &= g(w^{L_{l+1}} a^{L_l} + b^{L_{l+1}} + a^{L_l}) \end{aligned}$$

Now if  $w^{L_{l+1}} \neq b^{L_{l+1}}$  would become 0 then we would have

$$a^{L_{l+1}} = g(a^{L_l}) = a^{L_l}$$

This is because  $a^{L_l}$  is the  $\sigma$  ReLU on a  $\sigma$  the number is that number itself.

as we do

$$C^{(l+2)} = g(Z^{(l+2)} + C^{(l)})$$

thus both should have same dimensions.

Hence many times for result we use "Same" convolution.

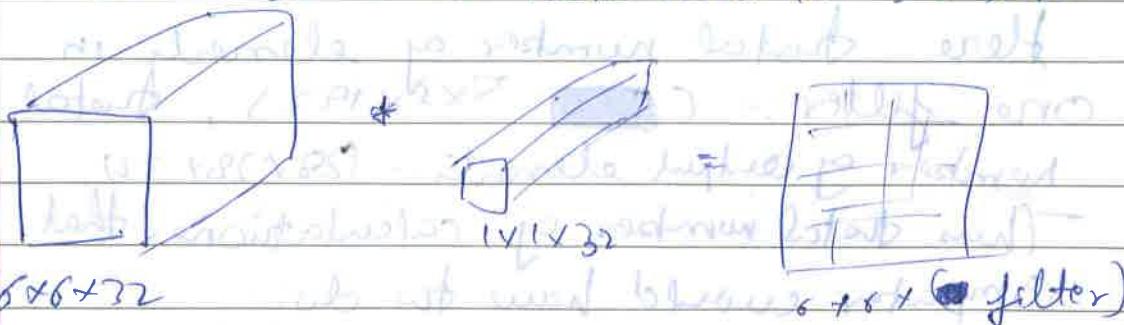
If they are not of same dimensions then we do:  
 $w^{(l+1)}$  where  $w$  is a matrix such that ~~dimensions~~  $w^{(l+1)}$  dimensions would become equal to  $Z^{(l+2)}$

→ Networks in networks and  $1 \times 1$  convolutions

$$\Rightarrow \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} * \begin{bmatrix} 2 \end{bmatrix} = \begin{bmatrix} 2 & 4 \\ 6 & 8 \end{bmatrix}$$

new it is not making much sense here (in 2-D).  
~~dimensions~~ because we are just multiplying each and every element by 2.

But it will make sense in 3-D like this:-

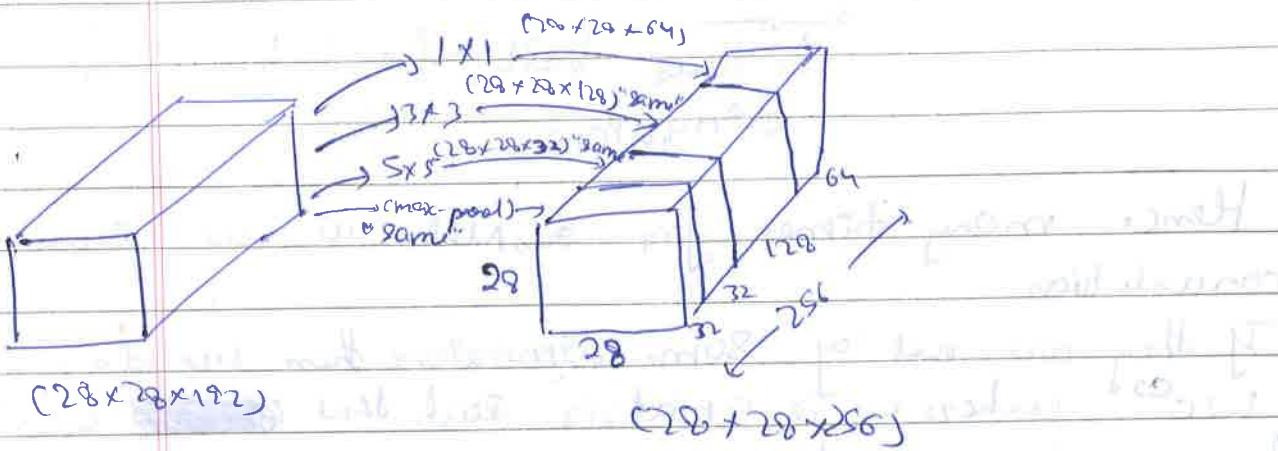


This idea of  $1 \times 1$  convolution is called network in network.

→ Using  $1 \times 1$  convolutions.

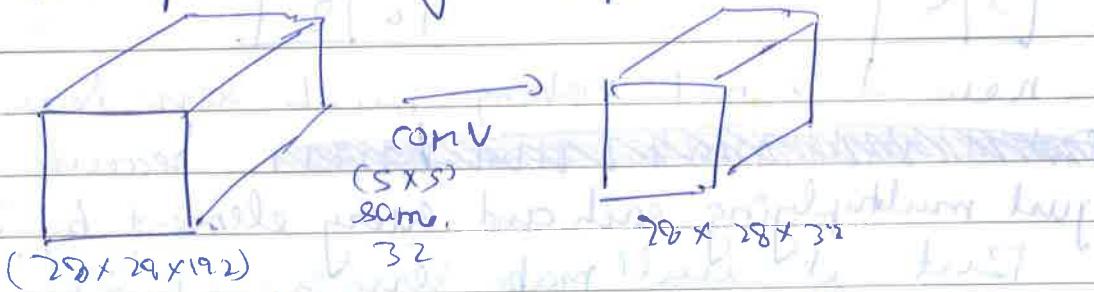
If you want to convert  $28 \times 28 \times 192$  to  $28 \times 28 \times 32$  then you can use 32  $(1 \times 1)$  filters.

→ Inception network motivation.



Here we are using different types of filters and then concatenating them.

→ The problem of computational cost.



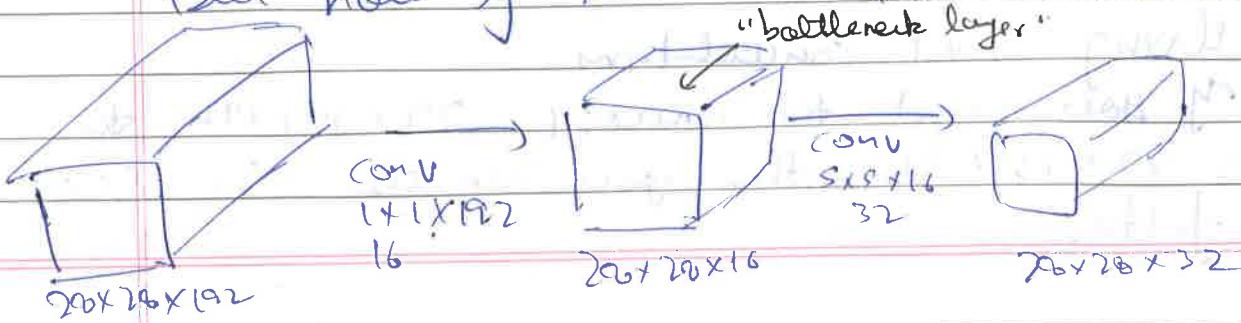
Here total number of elements in one filter :-  $(5 \times 5 \times 192)$ , total

number of output elements :-  $(28 \times 28 \times 32)$ .

Thus total number of calculations that computer would have to do:

$$28 \times 28 \times 32 \times 5 \times 5 \times 192 \approx 128 \text{ million}$$

But now if we use  $1 \times 1$  convolution.



So here the total number of computations would be:

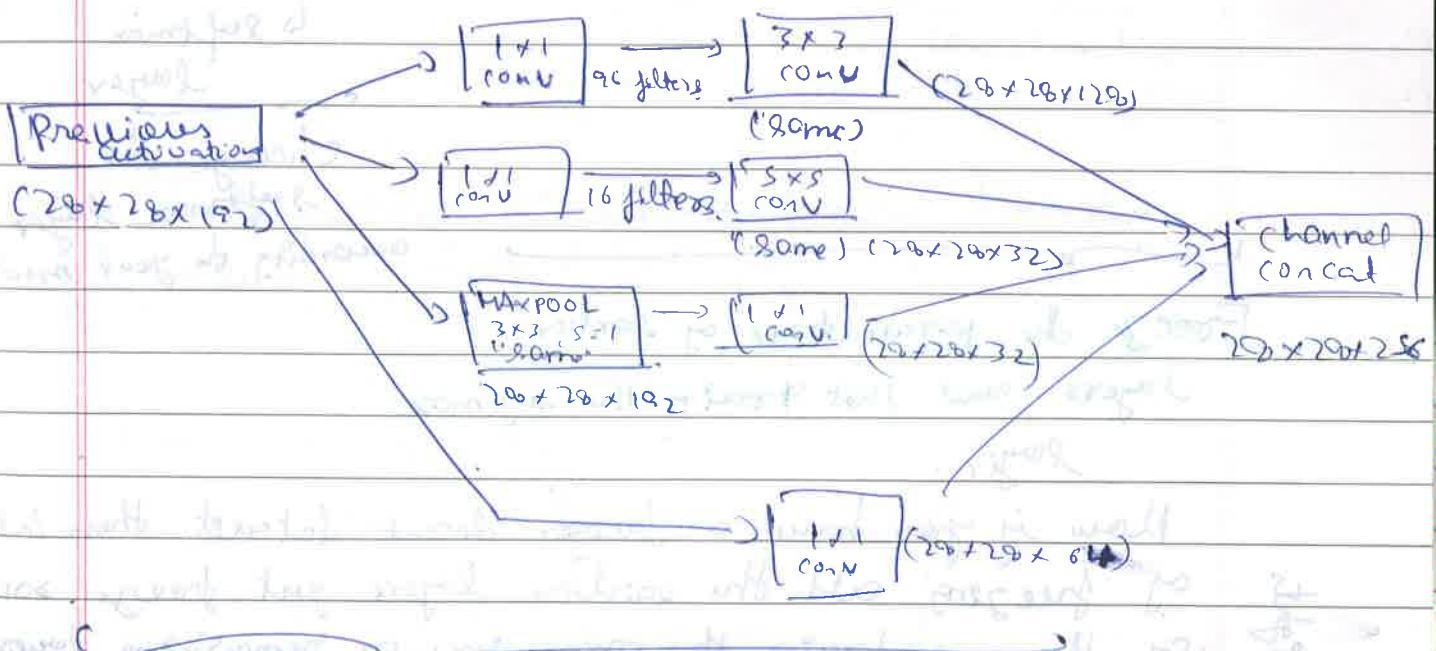
$$(28 \times 28 \times 16 \times 192) + (28 \times 28 \times 32 \times 8 \times 5 \times 16)$$

↓ 9.4 million      ↓ 10.0 million

↓ 19.4 million

Thus the computation cost is reduced by around 10 times.

→ Inception network



This is an inception block

→ Open source implementation of Comet  
You could use Resnet already available on GitHub.

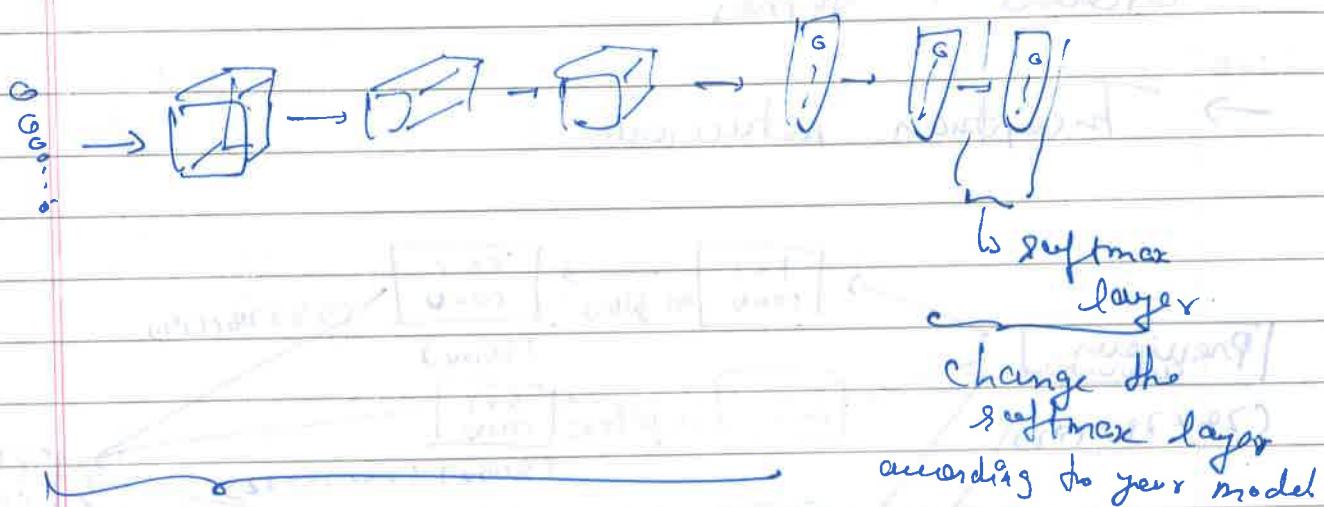
## → Transfer learning

Transfer learning is something that you must almost always do.

Let's say that you want to build a network that would classify a cat as → Tigger (Cat 1) and you don't have a very large dataset.

→ Misty (Cat 2)  
→ Neither

So you could use an open source trained neural network:



Freeze the parameters of earlier layers and just train the softmax layer.

Now if you have a ~~large~~ decent dataset then instead of freezing all the earlier layers just freeze some of them & train the parameters of remaining layers along with the softmax layers. Or you could add your own layers instead of using some of the layers of open source network.

If you have large enough dataset then you can ~~train the entire network~~ on your dataset and not freeze any of the layers.

## → Data Augmentation.

- (i) mirroring (using mirror image of object)
- (ii) Random cropping
- (iii) Rotation
- (iv) Shearing 
- (v) local warping
- (vi) Color shifting (adding some distortion to the R G B values).

Eg.  $\begin{matrix} +20 & -20 & \times 20 \end{matrix}$

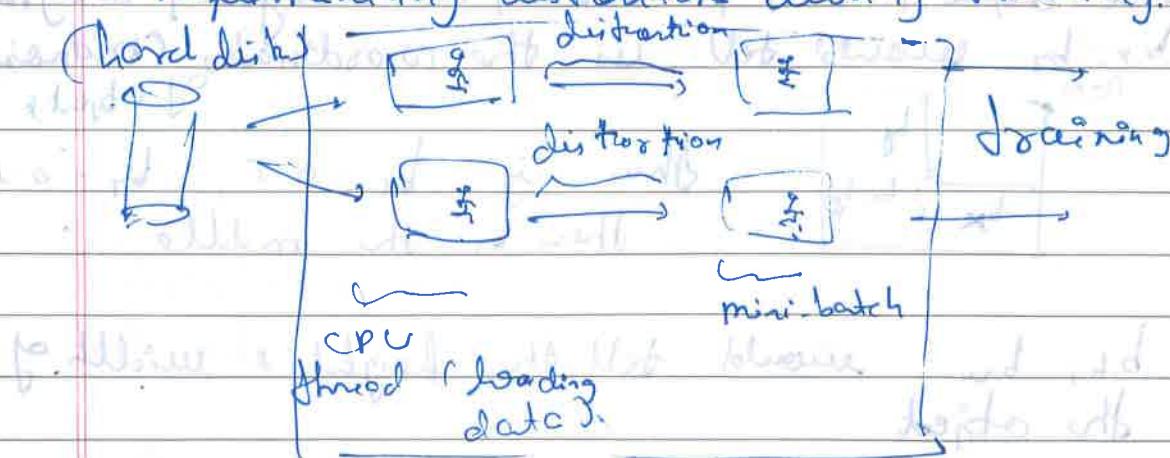
$\begin{matrix} R & G & B \end{matrix}$

→ This would make the image more people.

To immune the model to change in illumination we could use PCA ~~PCA~~ colour augmentation.

→ You could search online when need to use color shifting.

## (vii) Implementing distortions during training.



→ Ensembling.

→ Train several networks independently and average their outputs.

→ Multi-class at test time

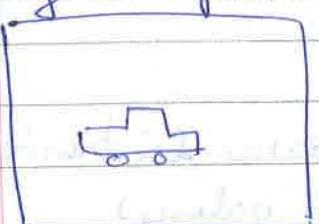
→ Run classifier on multiple versions of test images & average results.

→ This is used in competitions but not in production systems for real use.

# Specialization 4 week 3

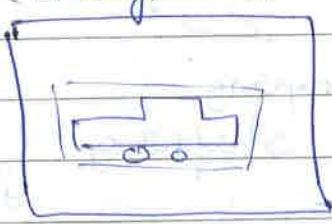
## → Object localization

(Image classification)



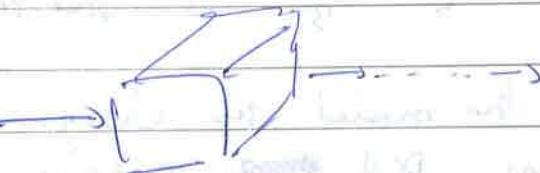
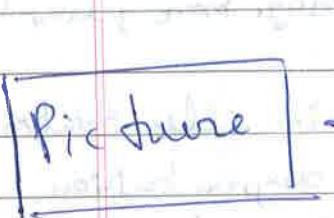
"Car"

(Classification and localization)



"Car"

Then we draw a boundary around the object also)



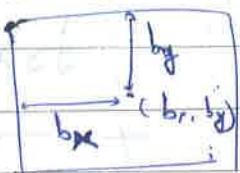
softmax

$b_x, b_y, b_h, b_w$

(boundary box)

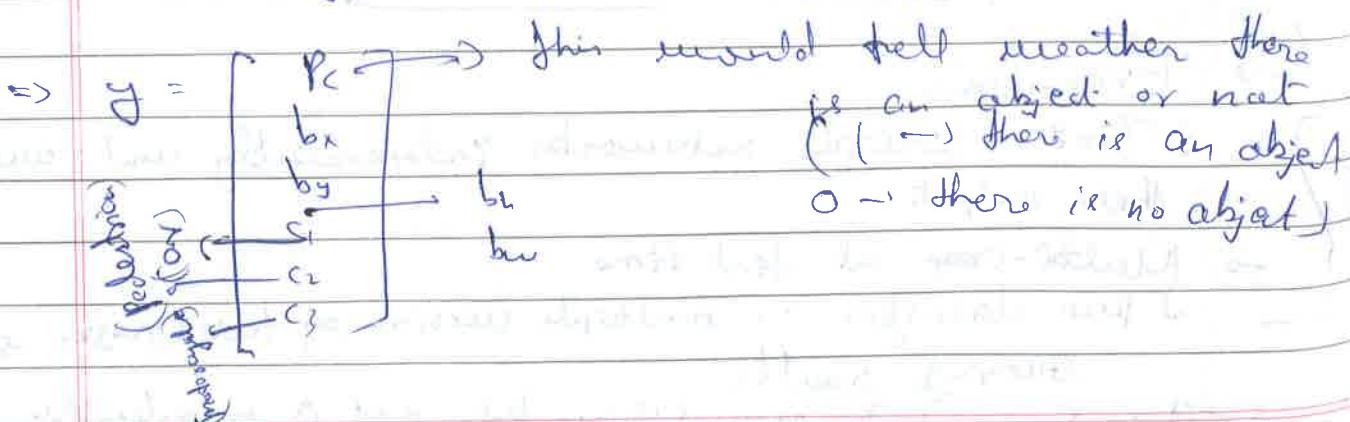
⇒ Here softmax layer would tell us the classification, the object (Car / pedestrian / motorcycle / background)

⇒  $b_x, b_y$  would tell us the co-ordinates (midpoint of object)



Then if  $b_x = 0.5, b_y = 0.5$ , then in the middle.

⇒  $b_h, b_w$  would tell the height & width of the object.



Eg. when there is a car:

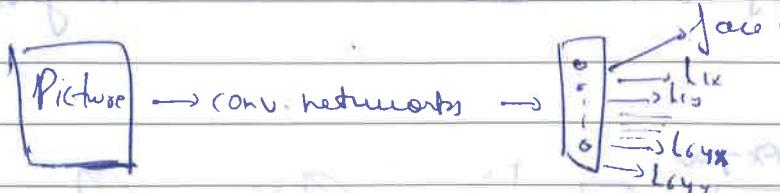
$$y = \begin{cases} 1 & \rightarrow y_1 \\ b_2 & \rightarrow y_2 \\ b_3 & \rightarrow y_3 \\ b_4 & \rightarrow y_4 \\ b_5 & \rightarrow y_5 \\ b_6 & \rightarrow y_6 \\ b_7 & \rightarrow y_7 \\ b_8 & \rightarrow y_8 \end{cases} \quad l(\hat{y}, y) = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2$$

Eg. when there is no object:

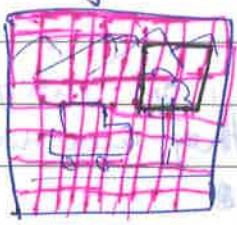
$$y = \begin{cases} 0 & \rightarrow \hat{y}_1 \\ ? & \rightarrow \hat{y}_2 \\ ? & \rightarrow \hat{y}_3 \\ ? & \rightarrow \hat{y}_4 \\ ? & \rightarrow \hat{y}_5 \\ ? & \rightarrow \hat{y}_6 \\ ? & \rightarrow \hat{y}_7 \\ ? & \rightarrow \hat{y}_8 \end{cases} \quad l(\hat{y}, y) = (\hat{y}_1 - y_1)^2$$

→ landmark detection:

For example we are making a face detection neural network. And for some reason we want to detect 64 points (landmarks) on the face. (Eg like the edges of face, edges of eyes, etc).



→ Object detection:



Suppose we have to detect whether there is a car and also show the boundaries of the car then we will take a neural network that given a cropped image of car tell whether it is a car or not.

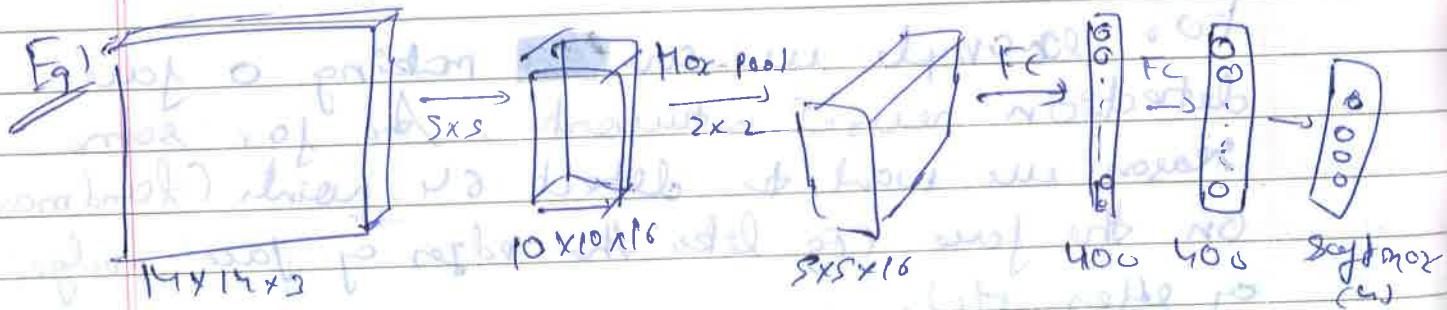
We crop out different part of the images one by one with a particular stride. And then pass into that neural network. Then it tells us that whether it is a car or not.

We can also change the size of the part that we were cropping out each time,

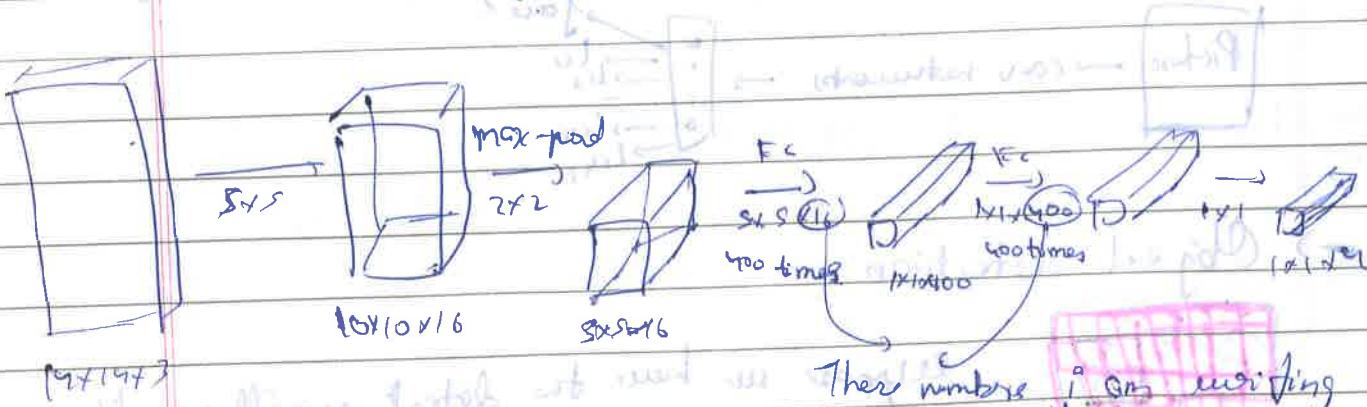
like initially we were choosing squares of  $(7 \times 7)$  only, now we will try squares of  $(3 \times 3)$  only.

Assume the image is of  $(20 \times 20)$  units.

→ Turning FC layers into convolutional layers.

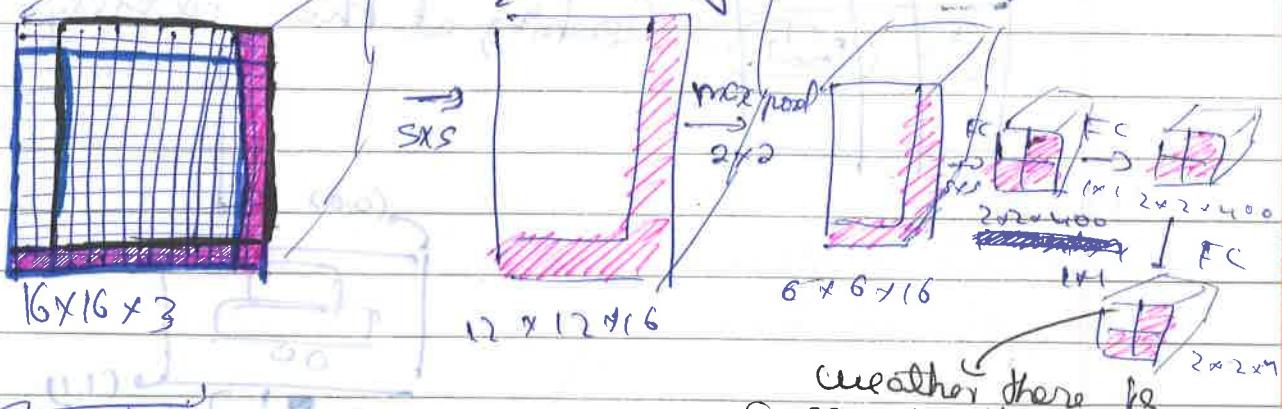


Now turning it into ~~conv~~ convnet layers.



These numbers are writing by myself, they were not in the video.

### Convolution implementation of sliding window.



Here we have divided the image into 4 windows.

whether there is a car in the first part or not.

The model tell us whether there is a car in the four windows (for each).

### Bounding box prediction

This object would be assigned to the first (0,0) part not to the second (0,1) part because the center of this object is in the first part.

100x100 pixel picture

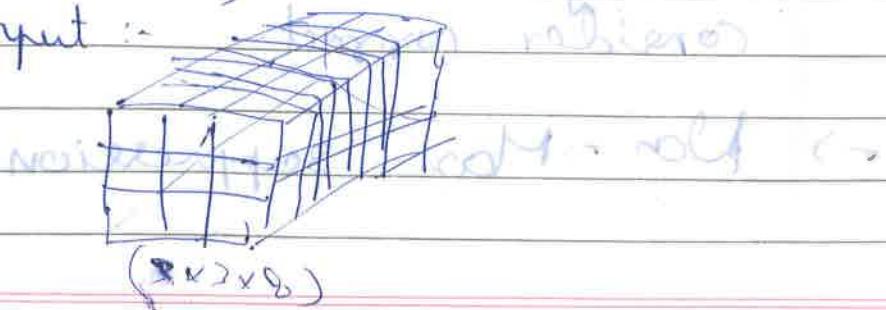
we divide into 9 parts.

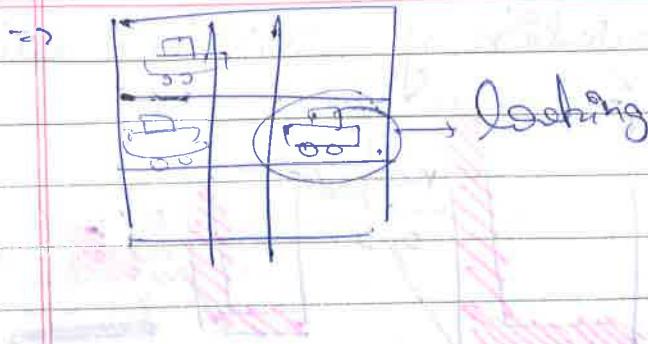
To sort on part we need to convert from

$$\begin{bmatrix} p_1 \\ b_1 \\ b_2 \\ b_3 \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

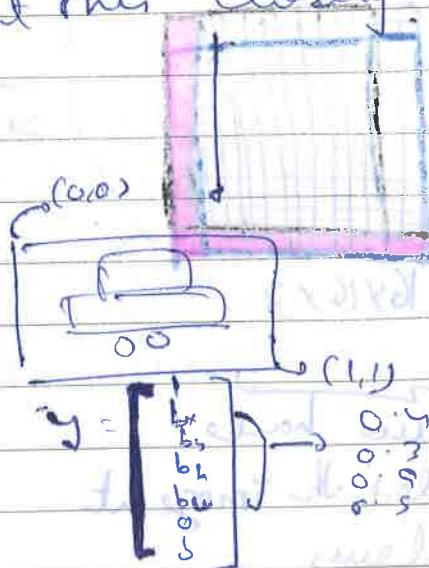
following we need to sort it.

Then output =

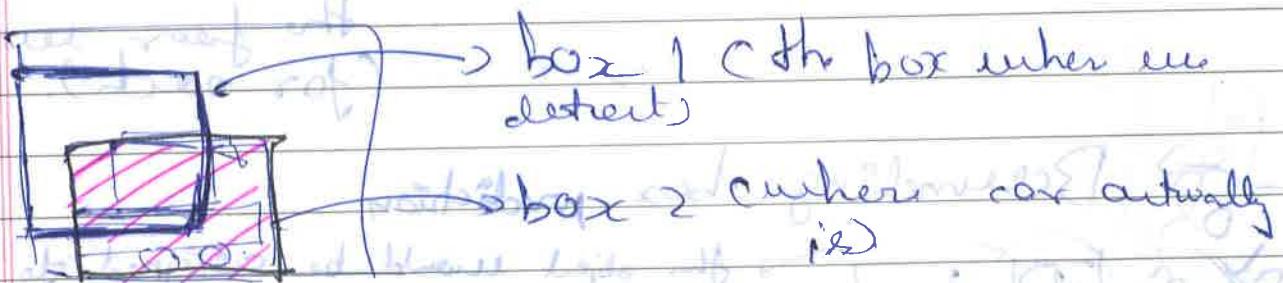




Looking at this closely



→ Intersection over union



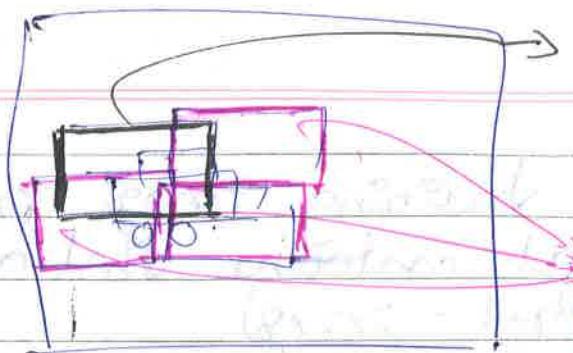
→ Intersection over union (IOU) :-

$\text{IOU} = \frac{\text{Size of intersection between box 1 \& box 2}}{\text{Size of union of box 1 \& box 2}}$

many times we use IOU to evaluate object ~~localization~~ localization.

If  $\text{IOU} \geq 0.5$  then we usually consider correct.

→ Non-Max suppression.



PAGE NO. \_\_\_\_\_  
 DATE \_\_\_\_\_

This has maximum  $P_c$ . Let's say for this  $P_c$  is 0.8.

These remaining boxes which are very much overlapping with the selected black box should be discarded.

(Here they are not being overlapped very much)

Then what we do is:-

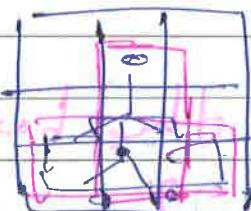
Discard all boxes with  $P_c \leq 0.6$

While there are any remaining boxes:-

→ Pick out the box with the largest  $P_c$  output as a prediction

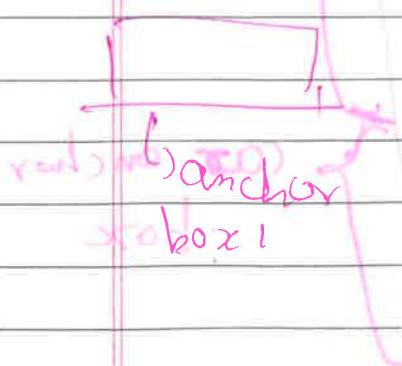
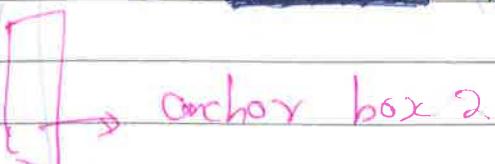
→ Discard any remaining box with  $IoU \geq 0.5$  with the box output in the previous step.

- Anchor boxes.
- overlapping objects.



midpoint  
of pedestrian  
and car

Here the midpoint of both the pedestrian and the car fall in the same anchor box.

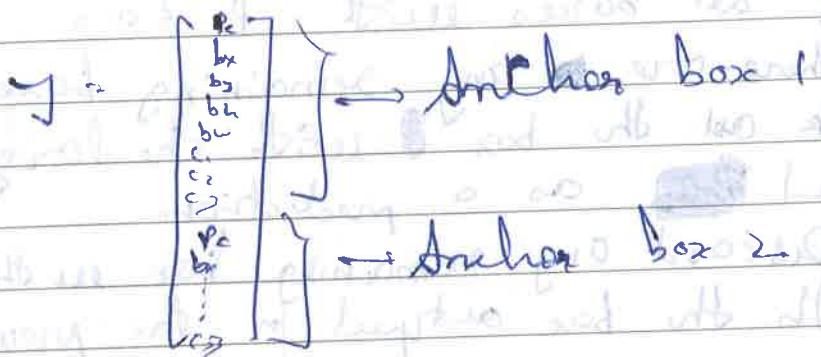


previously :-

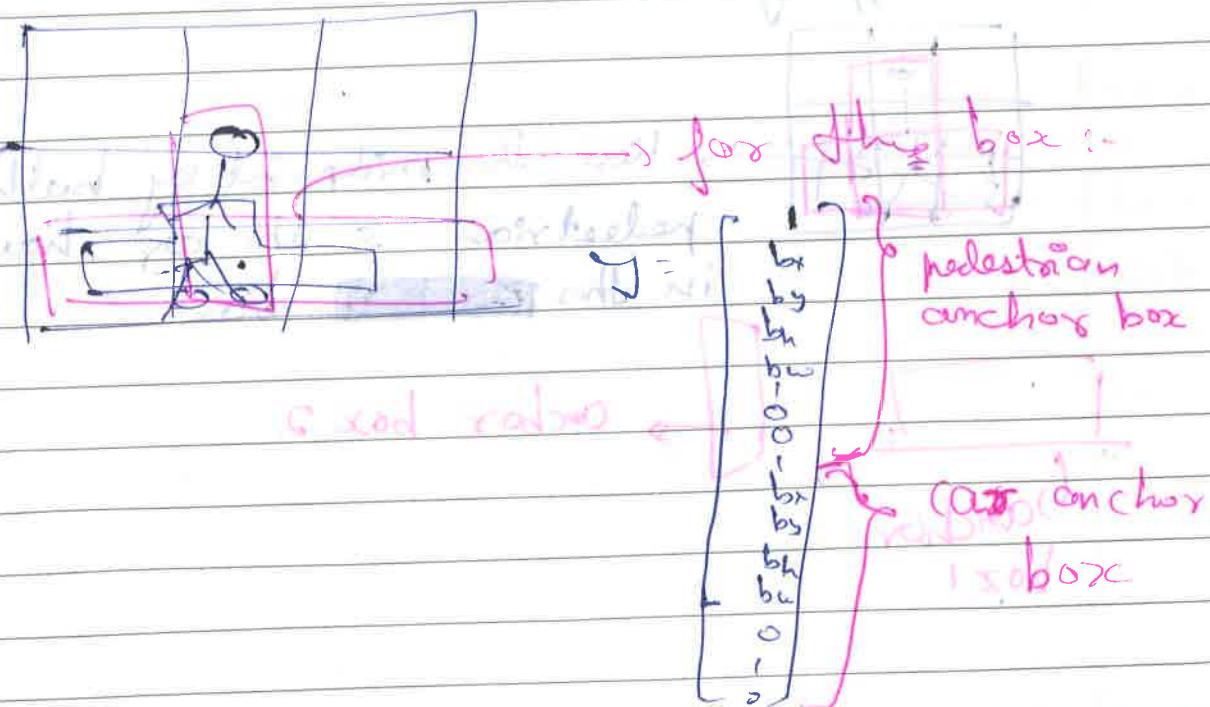
Each object in drawing image is assigned to grid cell that contains that object midpoint. (output :-  $3 \times 3 \times 8$ )

with two anchor boxes:-

Each object in drawing image is assigned to grid cell that contains object midpoint and anchor box for the grid cell with higher IoU. (output :-  $3 \times 3 \times 16$ )



then for this image :-



## Specialization 4. week 4.

### → Verification

1) Input image, name (ID)

2) Output whether the input image is that of the claimed person.

### → Recognition

1) Has a database of  $k$  persons.

2) Get an input image.

3) Output ID if the image is any of the  $k$  persons (or "not recognized")

### → One shot learning

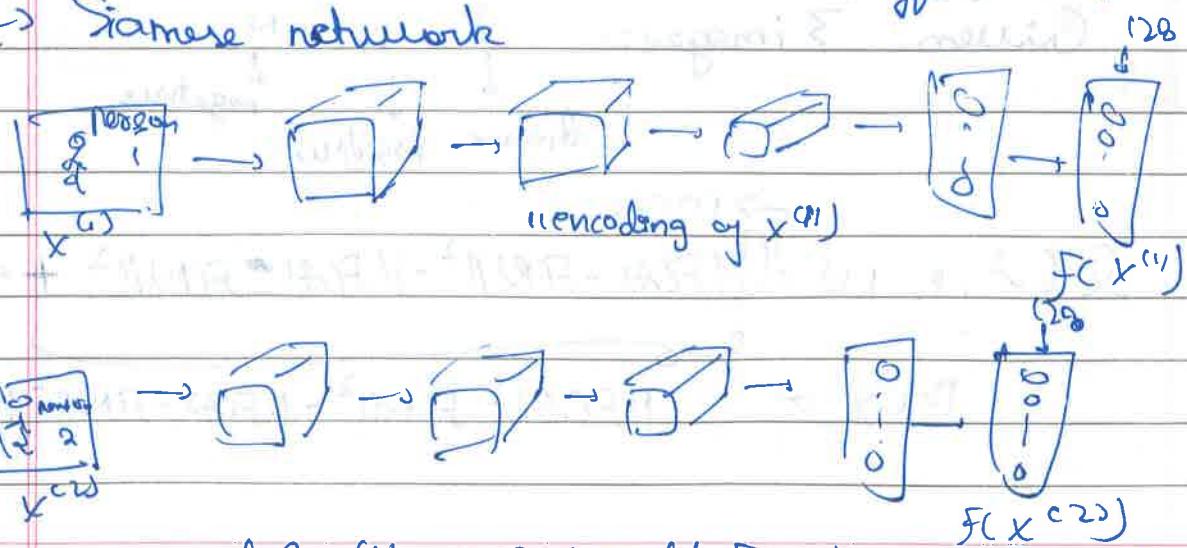
Many times we have just one or two images of a particular person. And we have to train our model to recognize that person within just that image.

So we use a "similarity" function.

$d(\text{img1}, \text{img2})$  = degree of difference between images.

If  $d(\text{img1}, \text{img2}) \leq \tau$  "same" } Verification  
 $\rightarrow \tau$  "different"

### → Siamese network



$$d(x^{c1}, x^{c2}) = \|f(x^{c1}) - f(x^{c2})\|_2^2$$

Parameters of  $\Phi$  define an encoding  $F(x^{(i)})$

learn parameters so that:

If  $x^{(i)}, x^{(j)}$  are the same person,

$\|F(x^{(i)}) - F(x^{(j)})\|^2$  is small

If  $x^{(i)}, x^{(j)}$  are different persons,

$\|F(x^{(i)}) - F(x^{(j)})\|^2$  is large

→ Triplet Loss

What we want is  $\rightarrow$  margin

$$\|F(A) - F(P)\|^2 + \alpha \leq \|F(A) - F(N)\|^2$$

$d(A, P)$   
Anchor image  
(initial image of a person)

$d(A, P)$   
positive image  
(image of the same person)

$d(A, N)$   
 $\downarrow$   
negative image  
(image of the another person)

$$\Rightarrow \|F(A) - F(P)\|^2 - \|F(A) - F(N)\|^2 + \alpha \leq 0$$

$\downarrow$

margin

⇒ Loss function

Given 3 images: -  $A, P, N$

$\downarrow$   $\downarrow$   $\downarrow$   
Anchor positive negative.

→ max

$$\mathcal{L}(A, P, N) = (\|F(A) - F(P)\|^2 - \|F(A) - F(N)\|^2 + \alpha, 0)$$

max of (  $\|F(A) - F(P)\|^2 - \|F(A) - F(N)\|^2 + \alpha, 0$  )

$$J = \sum_{i=1}^m l(A^{(i)}, p^{(i)}, N^{(i)})$$

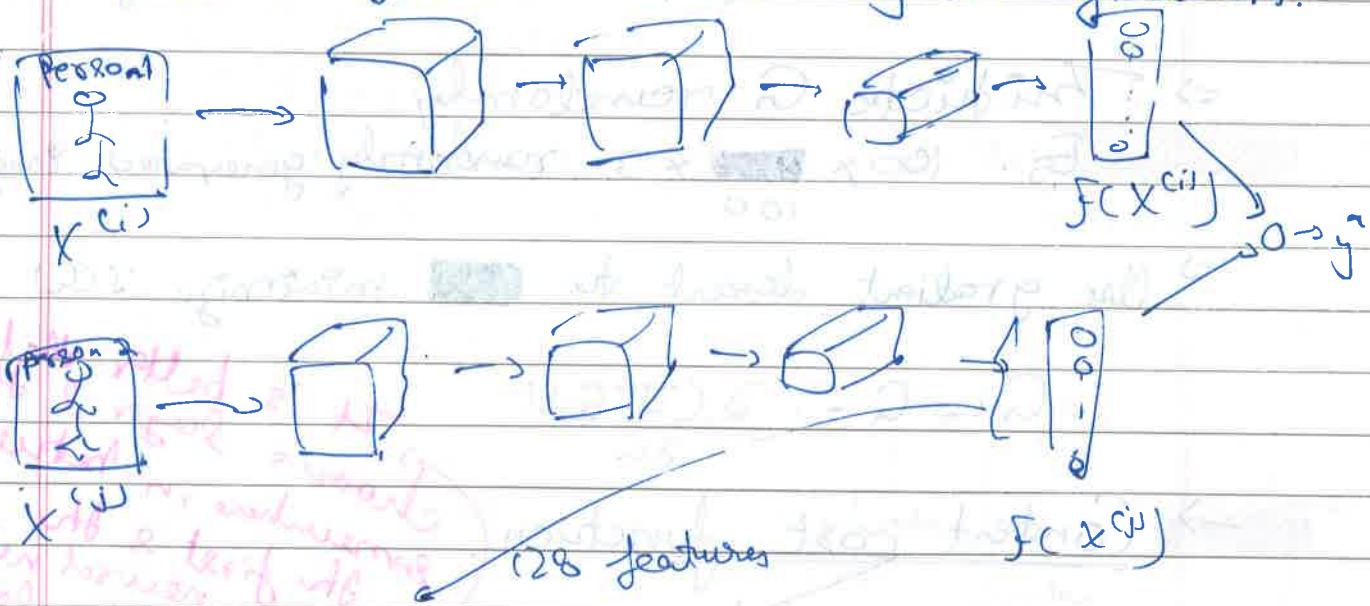
(cost function)

If you have training set of 10k picture  
eg 1k people then we have to select triplets  $(A, P, N)$

Now if we choose  $A, P, N$  randomly  
then  $d(A, P) + \alpha \leq d(A, N)$  is easily satisfied.

Then choose triplets that are "hard" to train on.

→ Face recognition and binary classification.



$$\hat{y} = \sigma \left( \sum_{k=1}^{128} w_k |f(x^{ci})_k - f(x^{sj})_k| + b \right)$$

or we could use

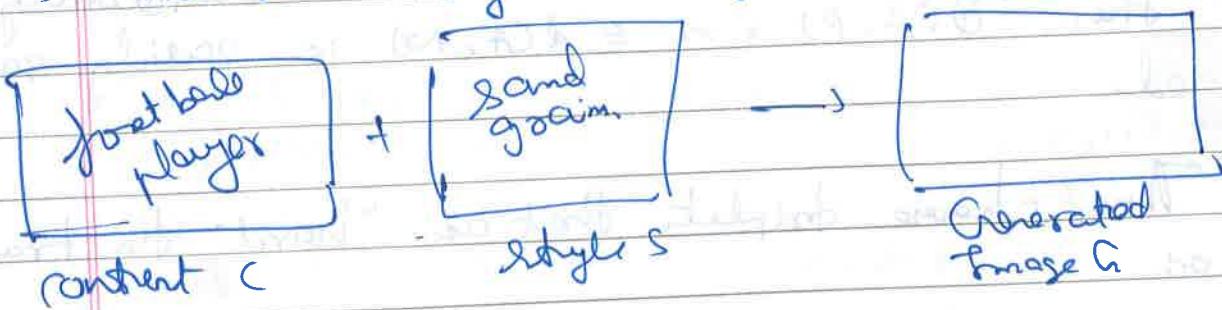
$$\left\{ \frac{(f(x^{ci})_k - f(x^{sj})_k)^2}{(f(x^{ci})_k + f(x^{sj})_k)} \right\}$$

→ So here we are using a pre-trained convolutional network and the building on top of it. ~~the idea of using a network trained on a different task and then applying it to a new task is called transfer learning which we are using here.~~

### → Neural style transfer

This adds style to content image. What onyma does in editing by footing. Like a background of sand on a football player would give a granular texture.

### → Neural style transfer cost function



$$J(G) = \alpha \cdot J_{\text{content}}(C, G) + \beta \cdot J_{\text{style}}(S, G)$$

→ Initialize G randomly.

Eg.  $100 \times \frac{100}{100} \times 3$  randomly generated image.

→ Use gradient descent to minimize  $J(G)$

$$G := G - \frac{\partial (J(G))}{\partial G}$$

### → Content cost function

$$J(G) = \alpha \cdot J_{\text{content}}(C, G) + \beta \cdot J_{\text{style}}(S, G)$$

• Say you use hidden layer  $l$  to compute content cost

• Use pre-trained Convnet (e.g. VGG network)

• Let  $\alpha_{\text{first}}$  and  $\alpha_{\text{second}}$  be the activation of layer  $l$  on the image.

It is better that yes  
choose layer l that is  
somewhere in between  
the first & the last  
neural network  
layers.

as edges and simple textures. The deeper layers tend to detect higher level features such as more complex textures and object classes etc. DATE

for the generated image "c" two have similar content as input image "c", then we should choose a layer somewhere in middle.

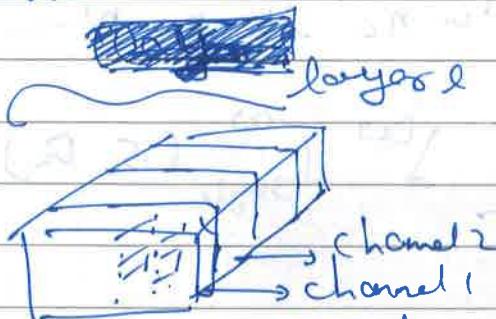
.) If  $a^{(0)}$  and  $a^{(L)}$  are similar, both images have similar content.

$$J_{\text{content}}(c, c) = \frac{1}{2} \|a^{(0)} - a^{(L)}\|^2$$

$$J_{\text{content}}(c) = \frac{1}{\text{number of all entries}} \sum (a^{(0)} - a^{(L)})^2$$

→ Style cost function.

Say that you are ~~using~~ using layer  $l$ 's activation to measure "style". Define style as correlation between activations across channels.



→ how well are these units in the ~~layer~~ channels are correlated to one another.

⇒ style matrix

let  $a_{i, j, k}^{(l)} = \text{activation at } (i, j, k)$ .

$i$  indicates height  
 $j$  indicates width  
 $k$  indicates channel

$C^{(l)}$  is  $n_h^{(l)} \times n_w^{(l)} \times n_c^{(l)}$

$$C^{(l)} = \sum_{i=1}^{n_h^{(l)}} \sum_{j=1}^{n_w^{(l)}} a_{i, j, k}^{(l)} a_{i, j, k}^{(l)}$$

for style image

→ for style image

In deeper layers of ConvNet, each channel corresponds to a feature detection. The weight matrix  $C_{ij}$  measures the degree to which the activations of different feature detectors in layer  $l$  vary or correlate together with each other.

$$G_{\text{reg}(G)} = \sum_{i=1}^{n_u} \sum_{j=1}^{n_w} G_{ijk} \left( \cdot \cdot G_{ijk} \right)$$

→ generated image

→ generated image.

This is "Gram matrix"

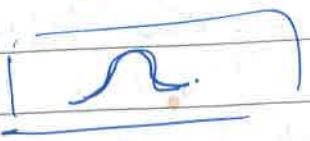
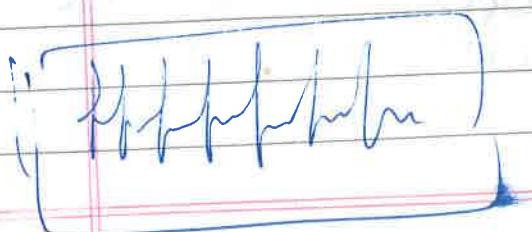
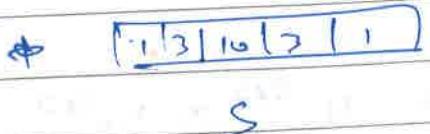
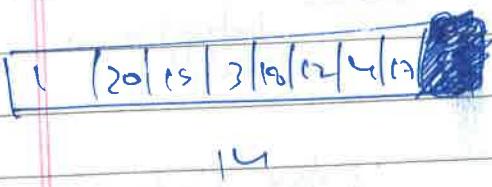
$$f_{\text{res}}^{\text{style}}(S, C_2) = \frac{1}{(2\pi f^{\text{res}}_r n^{\text{res}}_w n^{\text{res}}_c)} \cdot \sum_{R} \sum_{R'} (C_{R R'}^{\text{res}} - C_{R R'}^{\text{res}})$$

$$J_{style}(S, G) = \sum_{\alpha} J^{\alpha}_{style}(S, G)$$

$$J(G) = A \cdot J_{\text{content}}(G, G) + B \cdot J_{\text{style}}(G, G)$$

→ NB 3 3D generalization.

→ for 10 :-



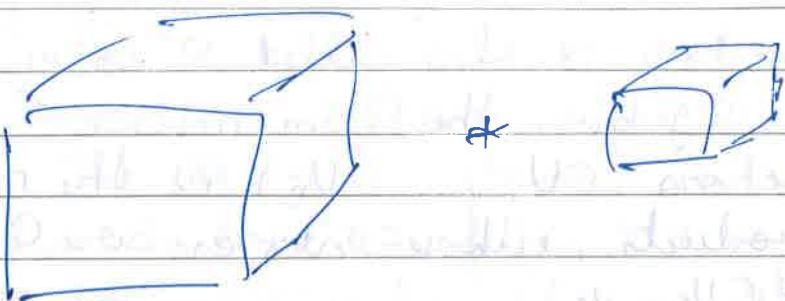
PAGE NO. \_\_\_\_\_ DATE \_\_\_\_\_

using ~~variables~~ filters

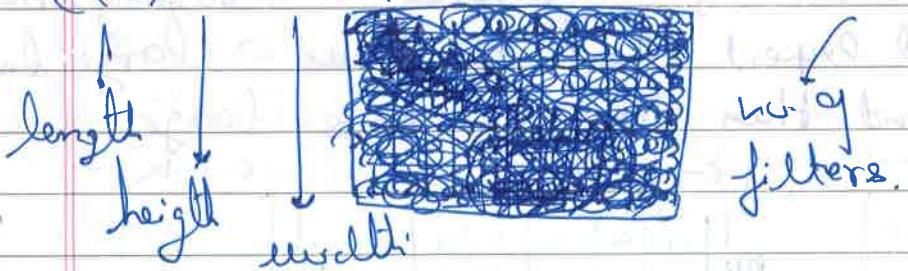
$$(8 \times 1) * (5 \times 1) = (4 \times 16) \quad *$$

$$(5 \times 16) \quad \boxed{C \times 16}$$

→ for 3-D



$$(4 \times 17 \times 17 \times 1) * (5 \times 5 \times 5 \times 16) = (10 \times 10 \times 10 \times 16)$$



→ Understanding neural style transfer a little more

e.g. VGG network

In neural style transfer we already have the network with all its parameters. Here we update the pixels of the input image that we put into the network rather than the parameters of that network. Earlier as we used to tune hyperparameters that used to determine how the network parameters would be updated, here we tune the hyperparameters that would determine how the

pixel values would be updated.

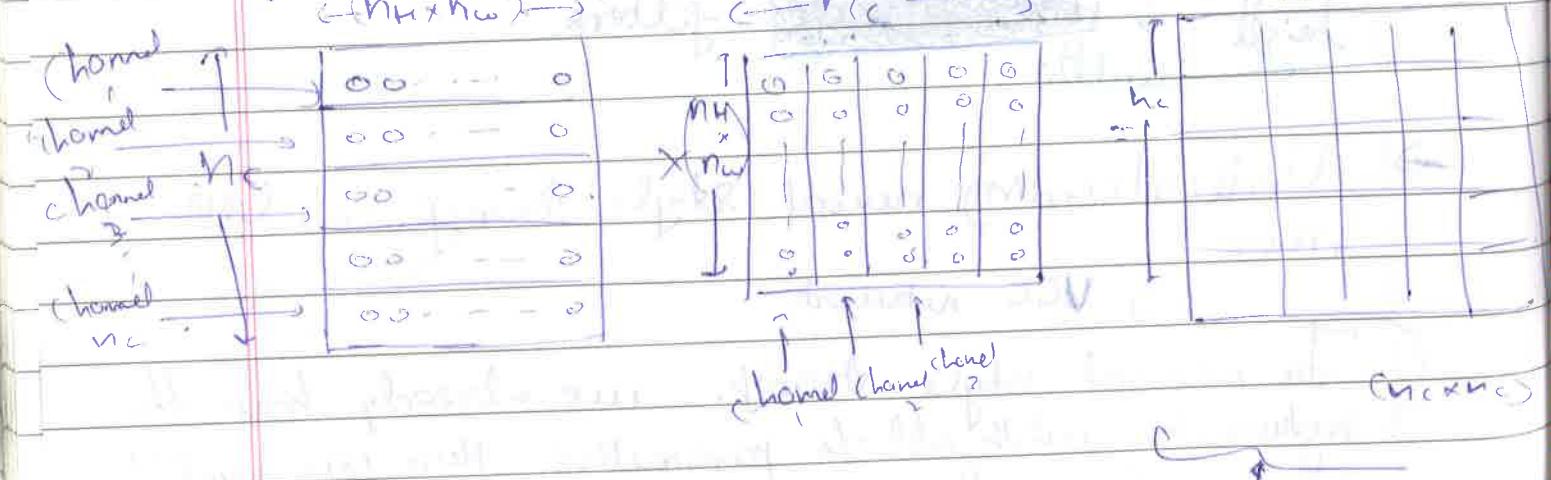
Content cost function ( $J_{content}(C, G)$ ):-

$$J_{content}(C, G) = \left( \frac{1}{4 \times n_H \times n_W \times n_C} \right) \sum_{\text{all } i, j} (C_{i,j} - G_{i,j})^2$$

all i, j entries

- ⇒ (i) The style matrix is also called a "Gram matrix".
- (ii) In linear algebra, the Gram matrix  $G$  of a set of vectors  $(V_1, \dots, V_n)$  is the matrix of dot products, whose entries are  $G_{i,j} = \underline{V_i^T V_j}$  =  $n_p \cdot \text{dot}(V_i, V_j)$ .

(iii) In other words,  $G_{i,j}$  compares how similar  $V_i$  is to  $V_j$ . If they are highly similar, you would expect them to have a large dot product, and thus for  $G_{i,j}$  to be large.



So in this matrix we would have all combinations of  $n_C \times n_C$  where  $n_i$  are all activations of a particular channel of a network and  $n_j$  are all activations of another particular channel of a network. G that

$$\begin{aligned}
 \text{np. lenalg. norm}(a) &= \sqrt{a_1^2 + a_2^2 + a_3^2 + a_4^2 + a_5^2 + a_6^2 + a_7^2 + a_8^2} \\
 &= \sqrt{(-4)^2 + (-2)^2 + (-2)^2 + (-1)^2 + (0)^2 + (1)^2 + (2)^2 + (3)^2 + (4)^2} \\
 &= \sqrt{7.245800} = 2.685800
 \end{aligned}$$

$G_{\text{gram}}_{i,j}$  : correlation

The Value  $G_{\text{gram}}_{i,j}$  measures how similar the activations of filter  $i$  are to the activations of filter  $j$ .

$\Rightarrow G_{\text{gram}}_{i,i}$  :- prevalence of patterns or textures.

$\Rightarrow$  The diagonal element  $G_{\text{gram}}_{i,i}$  measures how "active" a filter  $i$  is.

(ii) For example, suppose filter  $i$  is detecting vertical textures in the image. Then  $G_{\text{gram}}_{i,i}$  measures how common vertical textures are in the image as a whole.

(iii) If  $G_{\text{gram}}_{i,i}$  is large, this means that the image has a lot of vertical textures.

$$\Rightarrow \text{Style}_{\text{style}}(S, G) = \frac{1}{4 \cdot n_c^2 \cdot n_h^2 \cdot n_w^2} \sum_{i=1}^{n_c} \sum_{j=1}^{n_c} (G_{\text{gram}}^{(i)}_{i,j} - \bar{G}_{\text{gram}}^{(i)})^2$$

For a particular hidden layer  $l$  of the network.

$G_{\text{gram}}^{(i)}$  :- Gram matrix of "style" image

$G_{\text{gram}}^{(i)}$  :- Gram matrix of "generated" image.

$$\Rightarrow \text{Style}_{\text{style}}(S, G) = \sum_l \text{Style}_{\text{style}}(S, G)$$

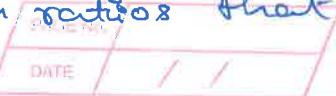
$$J(C) = \alpha \cdot J_{\text{content}}(C, C_0) + \beta \cdot J_{\text{style}}(C, C_0)$$

→ Image recognition.

We repeat the image verification technique with all the  $K$  images (of employees). We cross check the image of the face of the employee entering with the first image.

Compute that cost and assign that as min-cost. Now we cross check with 2nd image in the dataset. Compute cost, if the cost is less than min-cost then we would assign that cost to min-cost. We would repeat this for all the  $K$  pictures in database. Now if min-cost is less than the threshold value then we will allow that employee to pass.

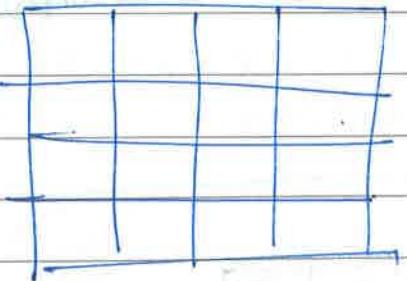
→ Anchor boxes are chosen by employing the training data to choose reasonable height/width ratios that represent the different classes.



→ Some more insight on anchor boxes :-

When using Anchor boxes, you can evaluate all object predictions at once. Anchor boxes ~~can~~ eliminate the need to scan an image with a sliding window that computes a separate prediction at every potential position. An object that uses anchor boxes can process an entire image at once, making real-time object detection systems possible.

The position of an anchor box is determined by mapping the location of the network output back to the input image. The process is replicated for every network output. The result produces a set of tiled anchor boxes across the entire image. Each anchor box represents a specific prediction of a class.



(4x4)

Input image.

Let's say that we <sup>define</sup> determine two anchor boxes. Then the output will have 2 anchor boxes per tile. Then we will have  $(2 \times 4 \times 4 = 32)$  anchor boxes in the output.

PAGE NO. \_\_\_\_\_  
DATE \_\_\_\_\_

before we eliminate them through YOLO's  
probability.

If the center / midpoint of an object falls into  
a grid cell, that grid is responsible for  
detecting that object

### → Understanding Anchor boxes with YOLO.

⇒ let's assume that we have an image with  
dimensions :- (m, 608, 608, 3)

⇒ Anchor boxes are chosen exploring the training  
data to choose reasonable height / width ratios  
that represent the different classes

⇒ For this problem, we will choose 5 anchor  
boxes (to cover the 80 classes).

⇒ The YOLO architecture is :-

Image (m, 608, 608, 3) → Deep CNN → Encoding (m,  $\frac{19}{1}, \frac{19}{1}, 5, 85$ )  
(example)(dimension)

{ 5 anchor  
boxes }

For each anchor box:-

$P_c | b_x | b_y | b_h | b_w | c_1 | c_2 | \dots | c_{80}$

Probability  
that there  
is an object

→ probability that the  
object is of class  $c_i$

thus Class scores :-

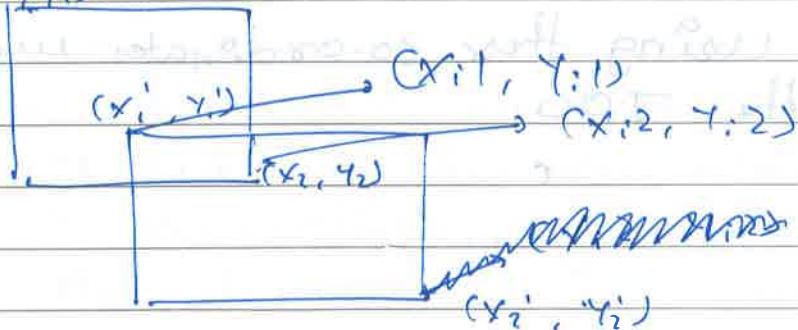
$$\text{scores} : P_c \cdot \begin{pmatrix} c_1 \\ c_2 \\ c_3 \\ \vdots \\ c_{80} \end{pmatrix} = \begin{pmatrix} P_c \cdot c_1 \\ P_c \cdot c_2 \\ P_c \cdot c_3 \\ \vdots \\ P_c \cdot c_{80} \end{pmatrix}$$

If we get  $P_c \cdot C_{30} = 0.6$ , then there is 60% probability that ~~an object~~ an object of class 30 exists, if we get  $P_c \cdot C_{33} = 0.73$ , then there is a probability of ~~73%~~ 73% that object of class 33 exists.

→ Then we will choose the object of maximum probability and we make a bounding box for that. Then in total we will have  $19 \times 19 \times 5 = 1805$  anchor boxes and then we will eliminate many of them ~~using~~ using non-max suppression.

→ Implementing IOU :-

$(x, y)$



Upper left corner of intersection box:-

~~$x_{11}$  :- maximum of  $x_1$  and  $x_2$  ;  $y_{11}$  :- minimum of  $y_1$  and  $y_2$  { as you go down the value of  $x$  increases}~~

Upper left corner of intersection box:-

(i)  $x_{11}$  :- rightmost of  $(x_1, x_2)$ , thus max of  $(x_1)$  and  $(x_2)$

(ii)  $y_{11}$  :- lowermost of  $(y_1, y_2)$ , thus min of  $(y_1)$  &  $(y_2)$  { as we go down, the value of  $y$  increases}

bottom right corner of intersection box:-

(i)  $x_{i2}$  :- leftmost of ( $x_2 \& x_{i'}$ ) ; thus the minimum of  $x_2 \& x_{i'}$ .

(ii)  $y_{i2}$  :- uppermost of ( $y_2 \& y_{i'}$ ) ; thus the minimum of  $y_2 \& y_{i'}$ .

The height of the intersection box would be:-

$$\boxed{y_{i2} - y_{i1}}$$

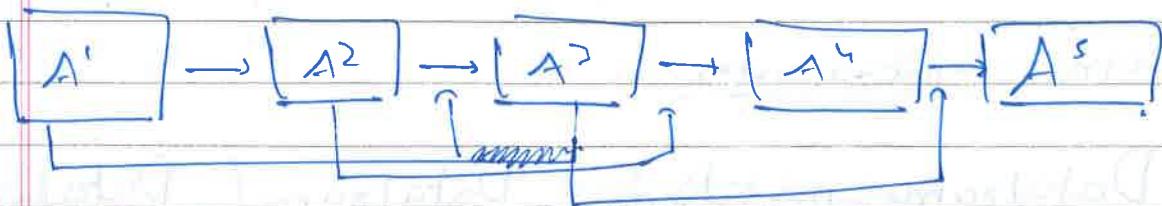
The length of the intersection box would be:-

$$x_{i2} - x_{i1}$$

If any of them is -ve then there is no intersection.

Now using these co-ordinates we would calculate the IOU.

→ Resnets



$$A^1 = g(\omega^1 A^0 + b^1)$$

$$A^2 = g(\omega^2 A^1 + b^2)$$

$$A^3 = g(\omega^3 A^2 + b^3) + A^1$$

$$A^4 = g(\omega^4 A^3 + b^4) + A^2$$

$$A^5 = g(\omega^5 A^4 + b^5) + A^3$$

$$\hookrightarrow g(\omega^5 [g(\omega^4 [g(\omega^3 [g(\omega^2 [g(\omega^1 A^0 + b^1) + A^1] + b^2) + A^2] + b^3) + A^3] + b^4) + A^4] + b^5) + g(\omega^3 A^2 + b^3) + A^1$$

So using Resnet we would be able to expand  $g$  in such way that this would keep into account the activations of previous layers such that it would not lead to vanishing or exploding gradients.

## Some notes ~~from~~ the API project

### ⇒ Some syntax notes:-

(i) `Dataframe-modified = Dataframe[ Dataframe["parameter"] == value ]`

This would return us a `dataframe` with all the values of the selected `parameter == value`

(ii) `df.drop(["col-1", "col-2", "col-3"], axis=1)`

This would drop `col-1`, `col-2` & `col-3` from the `dataframe`

(iii) `df.reset_index(drop=True, inplace=True)`

(iv) `df.dropna(axis=0, how='any', thresh=2)`

if `axis=0`, then it would drop rows and if `axis=1` then it would drop columns

if any NA value is present it would drop for example if `thresh=2` and if it were `how='all'` then all the values should be NA for it to drop then that row/column would be dropped if there were atleast 2 NA values in it

⇒ `df.dropna(subset=["col-1", "col-2"], inplace=True)`

This would state that we are to look into these columns when we would be looking for NA values to drop

→ df.set\_index(list, inplace=True)

PAGE NO. / / /

DATE / / /

(v) df.isna(); df["Parameter"].isna()

(vi) df.notna() → opposite of isna()

(vii) df[["parameter"]].isna().sum()

→ This would give the total number of samples that have missing value for that parameter

(viii) df[df.parameter == value].count()

This would return the number of samples that has that parameter = value

(ix) df.sort\_values(by = ["param\_1", "param\_2"])

(x) df\_shuffled = df.sample(len(df))  
df\_shuffled.reset\_index(drop=True, inplace=True)

or df\_shuffled = df.sample(frac=1)

(xi) df.loc[3]; df.iloc[3]

(xii) pd.crosstab(df["model"], df["doors"])

| doors  | 3 | 4 | 5 |
|--------|---|---|---|
| model  |   |   |   |
| BMW    | 0 | 0 | 1 |
| Honda  | 0 | 3 | 2 |
| Nissan | 1 | 3 | 0 |

(xiii) df.groupby(["make"]).mean()

| make   | adometre(kms) | doors |
|--------|---------------|-------|
| BMW    | —             | —     |
| Honda  | —             | —     |
| Nissan | —             | —     |

(xiv)  $df["Price"] = df["Price"].str.replace(['[141', '1.'], ''])$ . astype(int)

(xv) Suppose we are dealing with a classification problem with 2 results  $\rightarrow$  (0 & 1). We need to find how many 1's & 0's are there in the output. So we do :-

$df["result"].value_counts()$

$df["result"].value_counts(normalize=True)$

↳ This would give the distributions in  $\boxed{\quad}$  percentage.

(xvi)  $df + S \rightarrow$  same  
 $df + z \rightarrow$  same.

(xvii) Mlt. style. use ("recombine")

(xviii) car\_names = pd.Series(["BMW", "WACON", "SEAT"])  
 car\_colours = pd.Series(["Blue", "Black", "White"])

$df = pd.DataFrame(\{"car_names": car_names, "car_colours": car_colours\})$

$df["ready to use"] = True$

$df["Quantity"] = 4$

$df.rename(columns = \{"fruits": "my fruits", inplace = True\})$

$df["car_names"] = df["car_names"].str.lower()$

df.columns = new\_column\_names  $\rightarrow$  this is a good way to change the names of columns in dataframe

from basic operations on DataFrame

or you could do  $df = pd.DataFrame(\{"Carname": car_names, "Car-colour": car_colours\})$

.) Poem (Sheet 1 & 2) to get info about that syntax  
.) Poem (Sheet 1 →) to complete the several

(2d)  $\text{dFC}^{("odometer")}$ . $\text{fillna}(\text{dFC}^{("odometer")})$ .mean()  
implies = (true)

```
(x4) 1 animals = pd.Series(["cat", "dog", "lion"],  
                           index=[0, 1, 2])
```

(Contd) You can convert a 2-D numpy array to Pandas data frame. Index  $\rightarrow$  column names would take values of 0, 1, 2, ... .

(addii) np. Ones  $((2, 3)) \rightarrow \boxed{\text{[1, 1, 1]}} \quad [1, 1, 1], [1, 1, 1]$   
 np. zeros  $((1, 3)) \rightarrow [0, 0, 0]$

( $\alpha$ ,  $\beta$ ) np. random. scudent ( $0, 10, (5, 3)$ )

↓  
Inclusion  
Exclusion  
→ dimension of  
agency  
determined

Or we could have written this:

```
np.random.randint(6, 10, size=(5, 3))
```

np. `random.random((5,3))` → between 0

np. random. rand (5, 3) ↓ array of dimensions (5, 3) which has

$\xrightleftharpoons{\text{O}}$   $\xrightarrow{\text{H}_2\text{O}}$   $\xrightarrow[\text{L}]{\text{L}}$

( $\alpha$ ) np. unique (random\_array) or random\_array.unique()  
→ This would return all the unique elements

(xvi) np. square (a) ; np. sum (a) 

up.  $\text{sqrot}(a)$  ; np.mod(0, 2)  $\rightarrow$  0.1. 2

np. exp(a) ; np. log(c) ; np.var(a) ;

hp. max(a) i np. var(a) ; np. mean(a) ;  
np. std(a)

timeit

→ use this for python  
dataets

PAGE NO.

DATE

- (Q10.ii) i. **sum (a)**  
ii. timeit np.sum(a)

→ use this for numpy datasets

(Q10.iii) np.argsort (array)

np.argsort (array) → sort indices

np.argmax (array)

np.argmax (array)

np.argmax (array, axis = 1)

→ returns the index with  
maximum argument in rows.

np.argmin (array, axis = 0)

→ returns the index with minimum  
argument in columns.

(Q10.iv) df.aggregate ("min")

→ This would return the minimum value of  
each column

df.aggregate ("min", axis = 1)

→ would return the minimum value of each  
row, by default axis = 0

df.aggregate ([ "min", "max", "mean" ])

→ This would return the min, max  
& mean values of all columns.

df.aggregate ( { "Column-1": [ "sum", "mean" ],  
"Column-2": [ "sum", "max" ] } )

(axis=0) Suppose that we have a data frame like:

|   | A | B |
|---|---|---|
| 0 | 4 | 9 |
| 1 | 4 | 9 |
| 2 | 4 | 9 |

Then we do :- df.apply(np.sqrt)  
so we get :-

|   | A   | B   |
|---|-----|-----|
| 0 | 2.0 | 3.0 |
| 1 | 2.0 | 3.0 |
| 2 | 2.0 | 3.0 |

(axis=1) df.apply (lambda x: [1, 2], axis=1)

|   | [1, 2] |
|---|--------|
| 0 | [1, 2] |
| 1 | [1, 2] |
| 2 | [1, 2] |

Now :-

df.apply (lambda x: [1, 2], axis=1, result\_type='expand')

|   | 1 | 2 |
|---|---|---|
| 0 | 1 | 2 |
| 1 | 1 | 2 |
| 2 | 1 | 2 |

(axis=0) df.replace (to\_replace=0, value=df.mean(),  
inplace=True)

(axis=1) ~~cols = df.columns.tolist()~~ cols = df.columns.tolist()

~~cols = cols[-1:] + cols[:-1]~~ cols = cols[-1:] + cols[:-1]

df = df[cols]

{P. 7.0}

(b) (iii) np. where ( condition , x, y )

parameters:-

( condition:-

where True yield x otherwise y)

( x :- the value which will be yielded  
if condition is true )

( y :- the value which will be yielded  
if condition is false )

Eg 1



a = np. array ([ 0 )

>>> a

array ([ 0, 1, 2, 3, 4, 5, 6, 7, 8, 9 ])

>> np. where ( a < 5, 5, 10 \* a )

array ([ 0, 1, 2, 3, 4, 50, 60, 70, 80, 90 ])

→ condition

→ if condition is true this will be yielded ( x )

→ if condition is not true then this will be yielded ( y )

If we just give the condition then only those elements would be returned which fulfill that condition.

Eg

>> np. where ( a < 5 )

array ([ 0, 1, 2, 3, 4 ])

( add iv ) Plotting :-

{ p - 7. 0 }

Fig, (ax1, ax2, ax3, ax4, ax5)

= plt.subplots(nrows = 5, ncols = 1,  
figsize = (16, 32))

ax1.plot (df.index, df.param\_1, color = "red",  
label = "Param1 Vs index", marker = "o")

ax1.plot (df.index, df.result, color = "blue",  
label = "Param2 Vs index", marker = "o")

ax1.set\_xticks (ticks = np.arange (0, 40, 1))

ax1.set\_yticks (ticks = np.arange (0, 1, 0.05))

del df['fill\_between'] (df.index, df.param\_1, df.param\_2, where = (df.param\_2 > df.param\_1), color = "red")

ax1.fill\_between (df.index, df.param\_1, df.param\_2, where = (df.param\_2 > df.param\_1), color = "red", interpolate = True, alpha = 0.25)

ax1.fill\_between (df.index, df.param\_1, df.param\_2, where = (df.param\_2 < df.param\_1), color = "blue", interpolate = True, alpha = 0.25)

ax1.axhline (y = 0.5, color = "yellow", label = "yellow line", linewidth = 1.6, linestyle = "--")

ax1.axvline (x = 0.5, color = "black", label = "black line", linewidth = 1.6, linestyle = "-")

ax1.legend()

Q21. `set_title("-----")`;

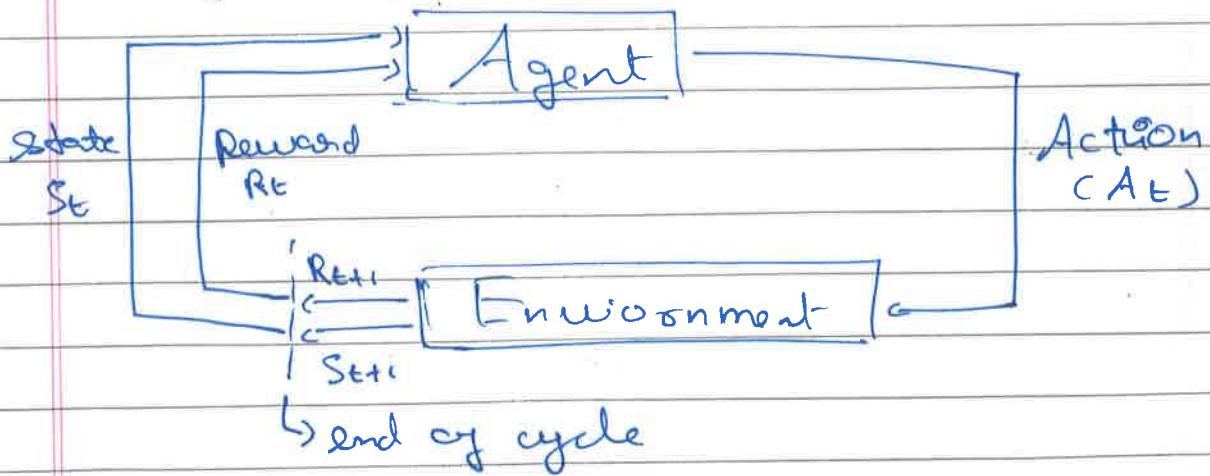
- Iteratively appending rows to a DataFrame can be more computationally intensive than a single concatenate. A better solution is to append those rows to a list and then concatenate the list with the original DataFrame all at once.
- `DataFrame.sample(frac = 0.5)`
-



## Reinforcement learning.

Reinforcement learning is an area of machine learning concerned with how software agents ought to take actions in an environment in order to maximize the notion of cumulative reward. Reinforcement learning is one of the three basic machine learning paradigms, alongside supervised learning and unsupervised learning.

Here we have a decision maker called an agent, that interacts with the environment it is placed in. These interactions occur sequentially over time. At each step, the agent will get some representation of the environment's state. Given this representation, the agent selects an action to take. The environment is then transitioned into a new state, and the agent is given a reward as a consequence of the previous action.



→ break

Let's break down the diagram into steps:

(i) At time  $t$ , the environment is in state  $S_t$

(ii) The agent observes the current state and selects action  $A_t$ .

(iii) The environment transitions to state  $S_{t+1}$  and grants the agent reward  $R_{t+1}$ .

(iv) The process then starts over for the next time step,  $t+1$ .

- Note,  $t+1$  is no longer in the future but is now the present. When we cross the dotted line on the bottom left, the diagram shows  $t+1$  transitioning into the current time step  $t$  so that  $S_{t+1}$  is now  $S_t$  and  $R_t$

This looping process was MDP (Markov decision process)

# R programming

1.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

2.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

3.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

4.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

5.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

6.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

7.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

8.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

9.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

10.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

11.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

12.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

13.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

14.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

15.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

16.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

17.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

18.  $\text{sum}(\text{c}(1, 2, 3, 4, 5))$  → 15

→ function as argument ( func(), a, b )  
 func (a, b) }

→ sentence ← function ( ... ) {

args = list ( ... )

a = ~~list~~ args [ 0 ]

b = args [ 1 ]

c = args [ 2 ]

print ( a, b, c, collapse = " " )  
 }

→ ~~unclasse~~ ( sys. time () )

this would return the number of seconds  
 passed since ~~list~~ 01-01-1970

→ unclasse ( sys. Date () )

number of days passed since 01-01-1970

→ unclasse ( as. Date ( " 1969-01-01 " ) )

↳ Ans: - 305

→ time is stored in two forms:-

↓  
 POSIXct

( stored as number of  
 seconds since 1970-01-01 )

↓  
 POSIXlt

( stored as a  
 list of seconds,  
 minutes, hours, etc )

→ gtr ( unclasse ( as. POSIXlt ( sys. time () ) ) )

↳ t ← as.POSIXlt ( sys. time () )

- To get a specific attribute  
 $t\$min$
- $weekday(t)$   
 ↳ would return today's DAY
- $month(t)$   
 ↳ would return today's month
- $quarters(t)$   
 ↳ would return the current quarter
- $t3 \leftarrow \text{"October 17, 1986 08:24"}$
- $t4 \leftarrow \text{strptime}(t3, "%Y-%m-%d %H:%M")$   
 ↳ There is a space here
- $\text{difftime}(\text{Sys.time}(), t4, \text{units} = \text{"days"})$
- $\text{head}(\text{data-frame})$
- $\text{dim}(\text{data-frame})$   
 ↳ "would return the number of rows and columns in data frame"
- $ls\_list \leftarrow \text{lapply}(\text{dataframe}, \text{class})$   
 ↳ list of input ↳ function  
 ↳ we want to apply on the list of inputs without parenthesis.
- (Creating Custom functions in R  
 Eg 1)

"% mult-add-one %"  $\leftarrow$  function (a, b) {  

$$(a + b) + 1 \}$$

4. 1. mult-add-on 1. 5. invert direction

21

- > viewing () (This would bring up documentation about the dataset.) (Just type viewing () and nothing else).
- > vapply () is a safer and faster choice for sapply. Faster because it does faster computation on larger datasets and safer because here we can explicitly mention in which format we want the output. And if we do not get in that format then an error will be thrown.

E.g.

vapply (flags, unique, numeric (1))  
 Here this should show error because we expect to get numeric output of vector length 1 but instead we get all unique values as output. So it throws error.

but if we do :-

vapply (flags, class, character (1))  
 This wouldn't show any error.

→ We can use `apply()` to split our data into groups based on value of some variable, and then apply some function to each group.

Eg.

`apply(flags$population, flags$red, mean)`

SS: 2 76.4

On an average 55.2 million people live in countries without red flag and 76.4 million people live in countries with red flag.

Eg/

`apply(flags$population, flags$landmarks, summary)`

For each and every landmarks the following will be returned:-

( min 1st Q. Median Mean 3rd Q. Max. )

→ When we apply :- `apply()` the result may not be very tidy so we apply ~~apply~~ ~~apply~~ `apply()` to simplify the results. `apply()` gives output as a list while ~~apply~~ ~~apply~~ applies the output to a vector, matrix or data frame.

Eg) `apply(fulsi_haldi, 2, c)`

|   |   |   |   |
|---|---|---|---|
| 0 | 0 | 1 | 2 |
| 4 | 5 | 7 | 9 |

min values

→ max values

→ class :- matrix

`apply(ingredients, 2, c)`

|       |     |
|-------|-----|
| fulsi | 0.4 |
| haldi | 0.5 |
| salt  | 1.7 |
| sugar | 2.9 |

class :- list  
dim :- null

→ `flag-values ← flags[11:17]`  
= to all bytes

→ `looping (unique_val, junction_element) element`

Here we are applying a custom function  
~~looping (unique\_val, junction\_element) element~~ on the data unique\_val,  
which return the second element of the  
data. So if data is a list of vectors  
so the will return the second element  
of each vector.

## Recurrent neural networks

Used in:- speech recognition, music generation, sentiment classification, DNA sequence analysis, Machine translation, Video entity recognition, Name entity recognition.

Eg1 Name recognition:

Harry Potter and Hermione Granger intended  
 Harry Potter and Hermione Granger intended  
 Harry Potter and Hermione Granger intended  
 Harry Potter and Hermione Granger intended

$\rightarrow y: 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0$

$y^{(1)} \ y^{(2)} \ y^{(3)} \ y^{(4)} \ \dots \ y^{(9)}$

$T_x = 9$  (Number of inputs)

$T_y = 9$  (Number of outputs)

### Vocabulary

|        |        |
|--------|--------|
| a      | 1      |
| aron   | 2      |
| and    | 3      |
| !      | 4      |
| Harry  | 4075   |
| ;      | 5      |
| Potter | 6830   |
| ;      | 6      |
| Zuley  | 10,000 |

Here  $x^{(i)}$  = Harry could be represented in hot-encoding as:-

|   |   |
|---|---|
| 0 | 1 |
| 0 | 1 |
| 0 | 1 |

$\rightarrow 4075$

|   |   |
|---|---|
| 0 | 1 |
| 0 | 1 |
| 0 | 1 |

$\rightarrow 6830$

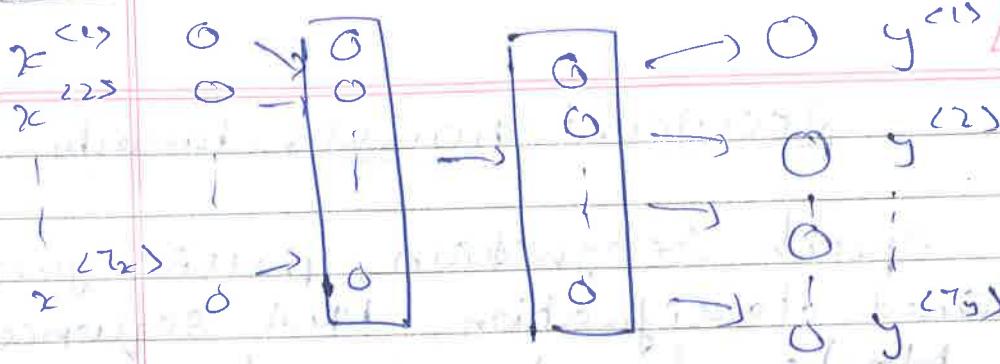
|   |   |
|---|---|
| 0 | 1 |
| 0 | 1 |

None-

hot-encoding

$\rightarrow$

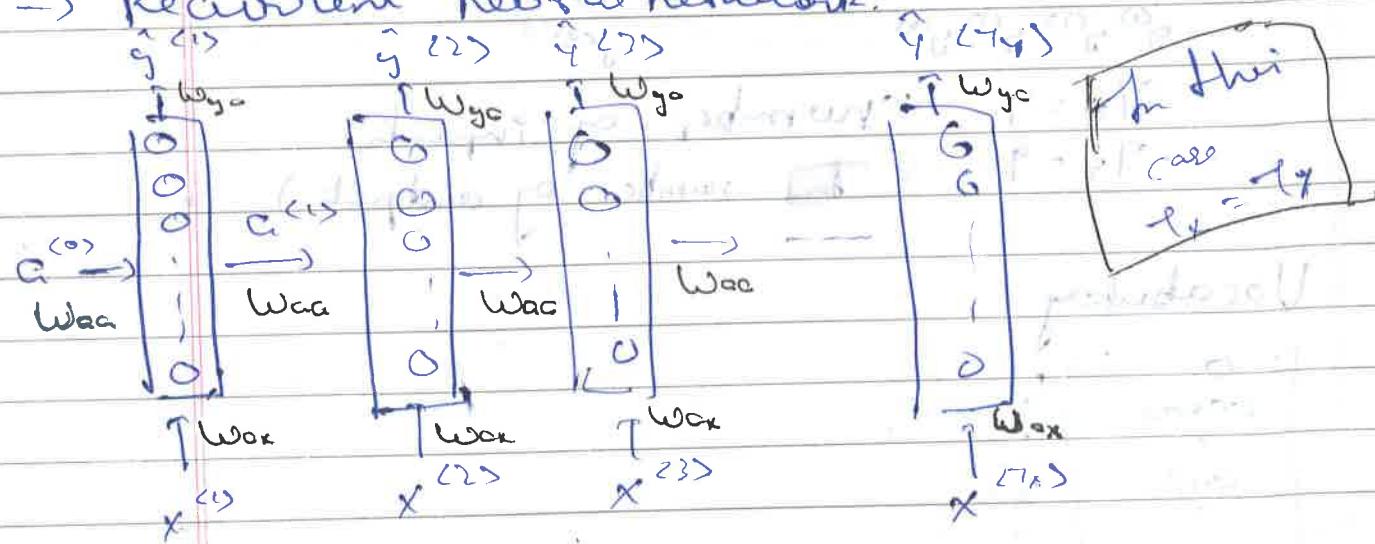
→ Ruby on standard network.



Probleme: ~~Was kann ich tun, wenn ich mich mit einem anderen~~

- (Inputs, outputs can be of different lengths in different examples)
  - (Doesn't share features learned across different positions of text)

→ Recurrent neural network.



→ Forward propagation

$$\vec{a}^{(0)} = \vec{0} \quad | \quad \vec{a}^{(1)} = g(C_{000}\vec{a}^{(0)} + C_{001}\vec{x}^{(0)} + \vec{b}_0)$$

$$\hat{y}^{(1)} = g(w_0a^{(1)} + b)$$

mostly we use Tanh or ReLU function as an activation function in RNN's

$$C^{(2E)} = g(C^{(2E-1)} + \text{loss} * x^{(E)} + b_0)$$

$$\hat{y}^{(t)} = g(w_0 e^{x^{(t)}} + b)$$

$$c^{(t+1)} = g(w_0 c^{(t+1)}, x^{(t+1)}) + b_0$$

■  $w_a = [w_{aa} \mid w_{ax}] \downarrow 100$   
 $\downarrow 100 \rightarrow \downarrow 10,000 \rightarrow$

$$= w_0 \downarrow 100,000,000$$

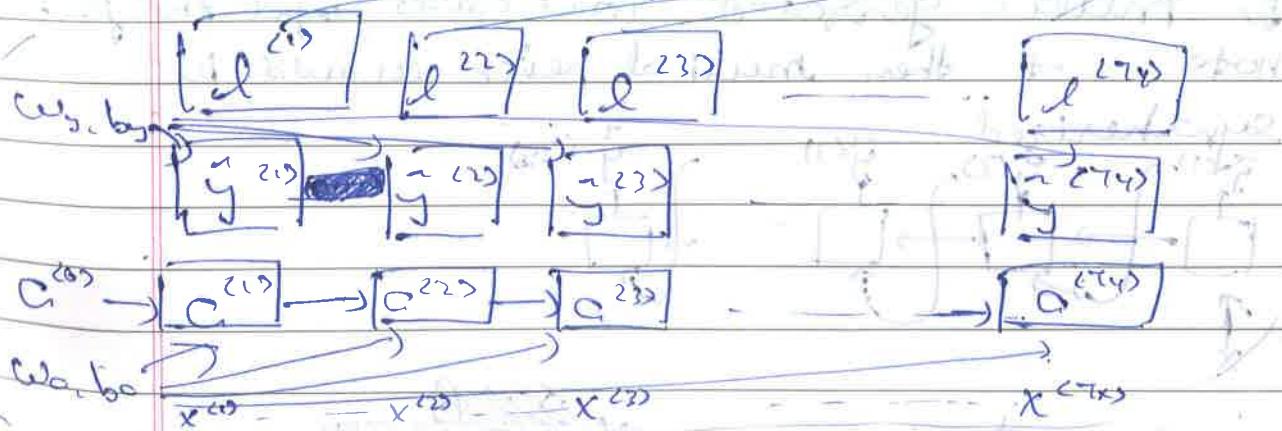
$$[c^{(t+1)}, x^{(t+1)}] = \begin{bmatrix} a^{(t+1)} \\ \otimes x^{(t+1)} \end{bmatrix} \downarrow 100 \downarrow 10,000 \downarrow 10,000 \downarrow 10,000$$

$$[w_{aa}, w_{ax}] \begin{bmatrix} a^{(t+1)} \\ x^{(t+1)} \end{bmatrix} = w_{aa} a^{(t+1)} + w_{ax} x^{(t+1)}$$

Thus forward propagation



Thus forward propagation



$$l^{(t+1)}(y^{(t)}, y^{(t)}) = -y^{(t)} \log y^{(t)} - (1-y^{(t)}) \log (1-y^{(t)})$$

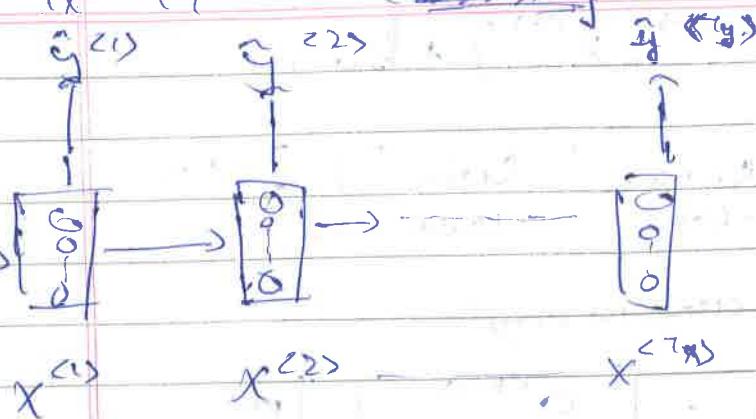
$$l(\hat{y}, y) = \sum_{t=1}^T l^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

AP-T-G

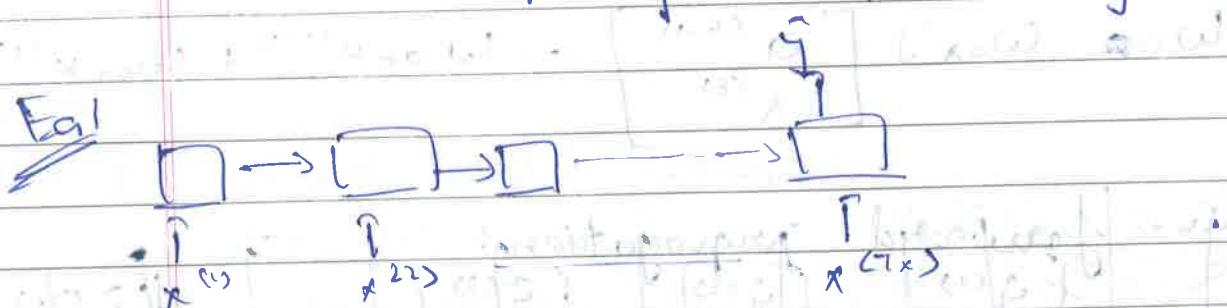
→ Example of RNN architectures.

|          |     |
|----------|-----|
| PAGE NO. | 111 |
| DATE     |     |

Eg)  $T_x = T_y$  (many-to-many)

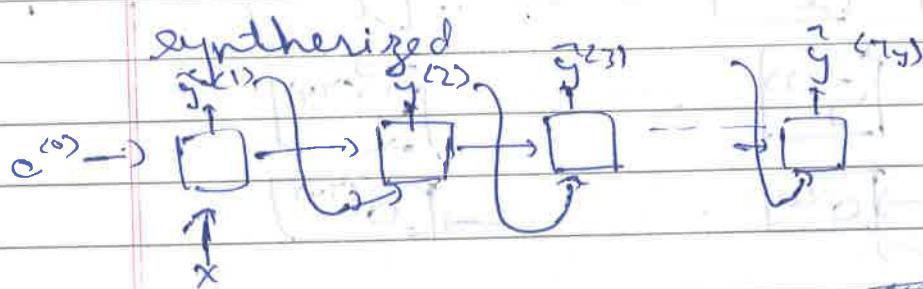


This is an example of many-to-many architecture.



This is a many-to-one architecture. Sentiment Analysis could be an example.

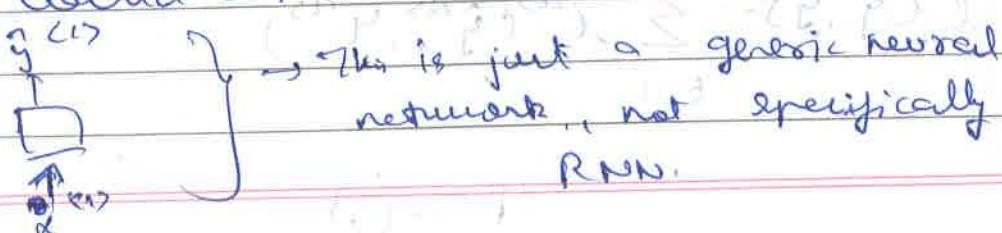
Eg) One-to-many architecture. Example could be music generated. You could input the first note and then musical piece would be



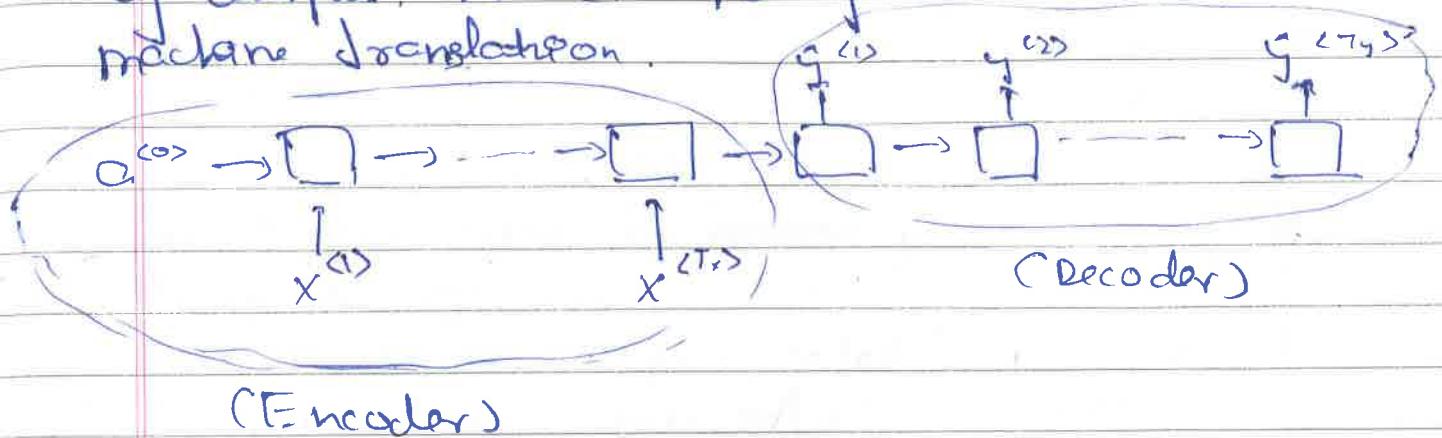
one-to-many.

~~where X is initial~~

Eg) we could also have one-to-one.



Eg) example of many-to-many architecture with the number of inputs not equal to the number of outputs. An example of this could be machine translation.



→ Autoencoders.

Unsupervised neural networks ~~.....~~

→ What is language modelling.

Speech recognition:

The apple and ~~pear~~ salad

The apple and pear salad.

Now ~~probabilistic~~ language model will tell which of the two should be the one when both of them sound exactly the same.

So it will return the ~~max~~ probability of the sentence being the right one.

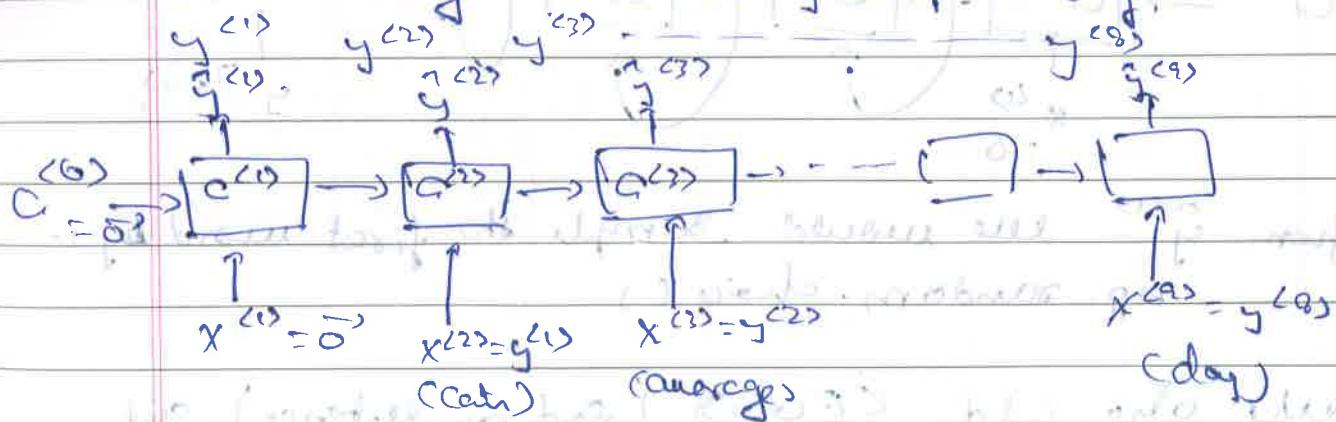
$$p(\text{The apple and } \cancel{\text{pear}} \text{ salad}) = 3.2 \times 10^{-13}$$

$$p(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$$

Training set :- large corpus of English text.

Tokenize.

(On average 15 hours of sleep ⋅ day.



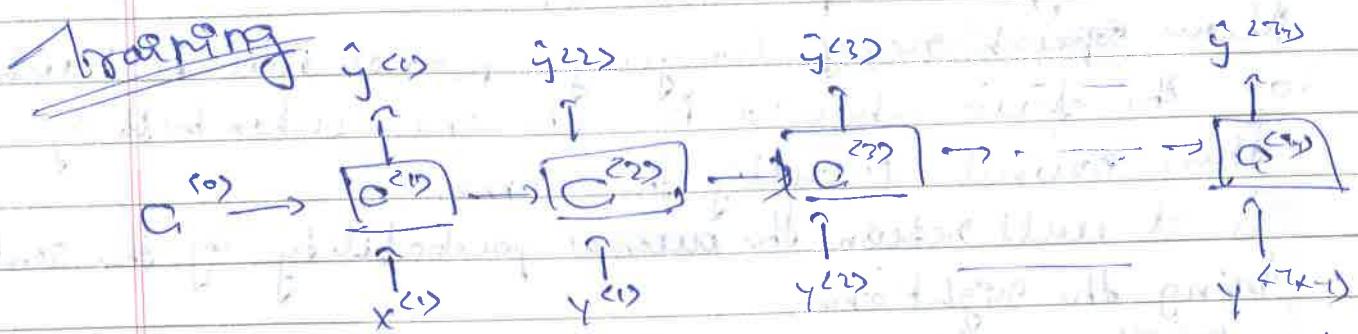
$\hat{y}^{(t)}$  would be softmax vector telling the probability of each and every word of the vocabulary. we add  $\langle \text{EOS} \rangle$  and  $\langle \text{UNK} \rangle$  token in our corpus.

Vocabulary also

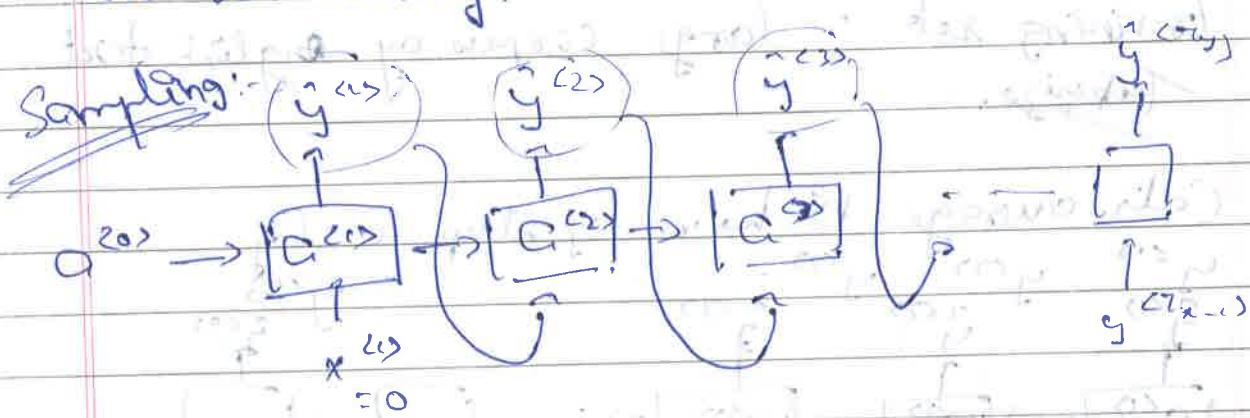
$$d(y^{(t)}, y^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$
$$d = \sum_i d(y_i^{(t)}, y^{(t)})$$

Then the probability that the given sentence is  $p(y^{(1)}, y^{(2)}, y^{(3)})$  is  $p(y^{(1)}) \cdot p(y^{(2)}|y^{(1)}) \cdot p(y^{(3)}|y^{(1)}, y^{(2)})$

# Sampling A Sequence from a Trained RNN.



here each and every  $y$  would be a software vector giving the probability of each and ~~every~~ word from the vocabulary.



from  $y^{(1)}$  we would sample the first word by :-  
np. random.choice()

We also add `<EOS>` {end of sentence} and `<UNK>` {unknown} token in our vocabulary too.

→ ~~word~~ Till now we were looking at word-level language model. Now let's look at character-level language model.

Vocabulary = [c, b, c, z, ., b, i, A, ..., Z]

Here each and every character is given as output.

→ Vanishing gradients with RNN's

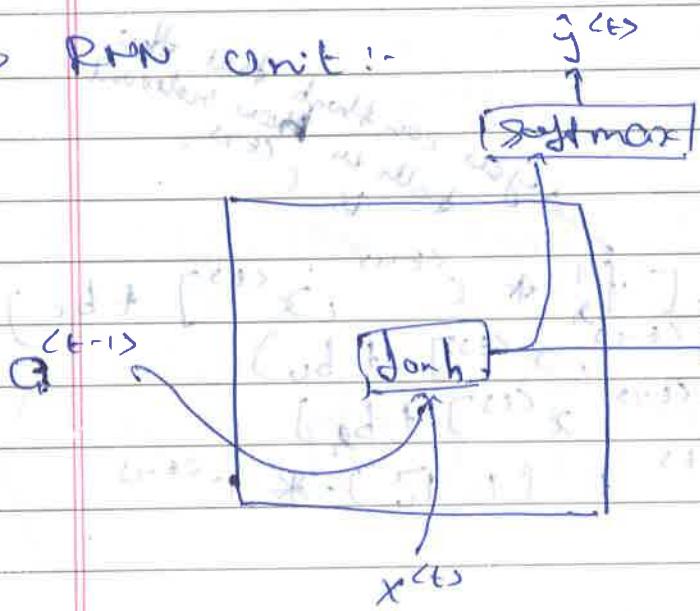
PAGE NO.

DATE

- Eg. The [cat] which [was] full.  
→ if singular cat then was  
The [cats] which [were] full  
→ if plural cat then were

These words much earlier in the sentence can have effect on words much farther in the sentence. Thus RNN must also have long range dependencies. Vanishing gradients is a big problem in RNN

→ RNN Unit:-



→ GRU (Simplified) (Gated Recurrent Unit)

(= memory cell. For this example see cell consider  $a^{t-1}$  and  $c^{t-1}$ )

$$\tilde{c}^{t-1} = \tanh(\omega_c [c^{t-1}, x^{t-1}] + b_c)$$

$$\Gamma_u = \sigma(\omega_u [c^{t-1}, x^{t-1}] + b_u)$$

↳ "update"

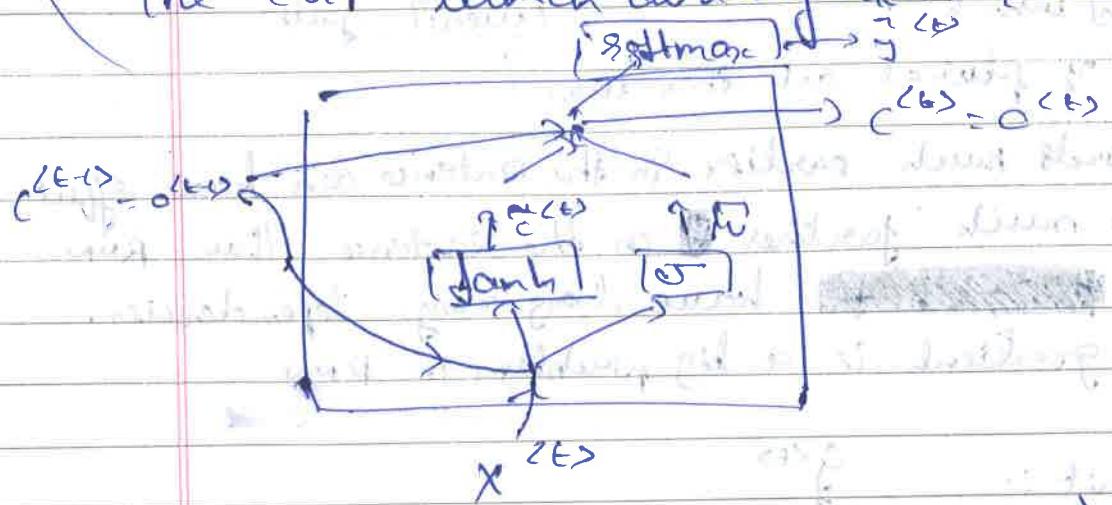
$$c^{t-1} = \Gamma_u * \tilde{c}^{t-1} + (1 - \Gamma_u) * c^{t-1}$$

element wise

$$\Gamma_U = 1 \quad \Gamma_U > 0 \quad \Gamma_U > 0 \dots \quad \Gamma_U = 1$$

$$\tilde{c}^{(t)} = 1$$

The cat, which already ate, was full.



$\Rightarrow$  Full CRU

$$\tilde{c}^{(t)} = \text{Janh}(\mathbf{w}_c [\Gamma_U, x^{(t)}] + b_c)$$

$$\Gamma_U = \sigma(\mathbf{w}_u [\mathbf{c}^{(t-1)}, x^{(t)}] + b_u)$$

$$\Gamma_F = \sigma(\mathbf{w}_f [\mathbf{c}^{(t-1)}, x^{(t)}] + b_f)$$

$$c^{(t)} = \Gamma_U * \tilde{c}^{(t)} + (1 - \Gamma_U) * c^{(t-1)}$$

$\rightarrow$  LSTM

$$\tilde{c}^{(t)} = \text{Janh}(\mathbf{w}_c [\mathbf{a}^{(t-1)}, x^{(t)}] + b_c)$$

~~update gate~~  $\Gamma_U = \sigma(\mathbf{w}_u [\mathbf{a}^{(t-1)}, x^{(t)}] + b_u)$

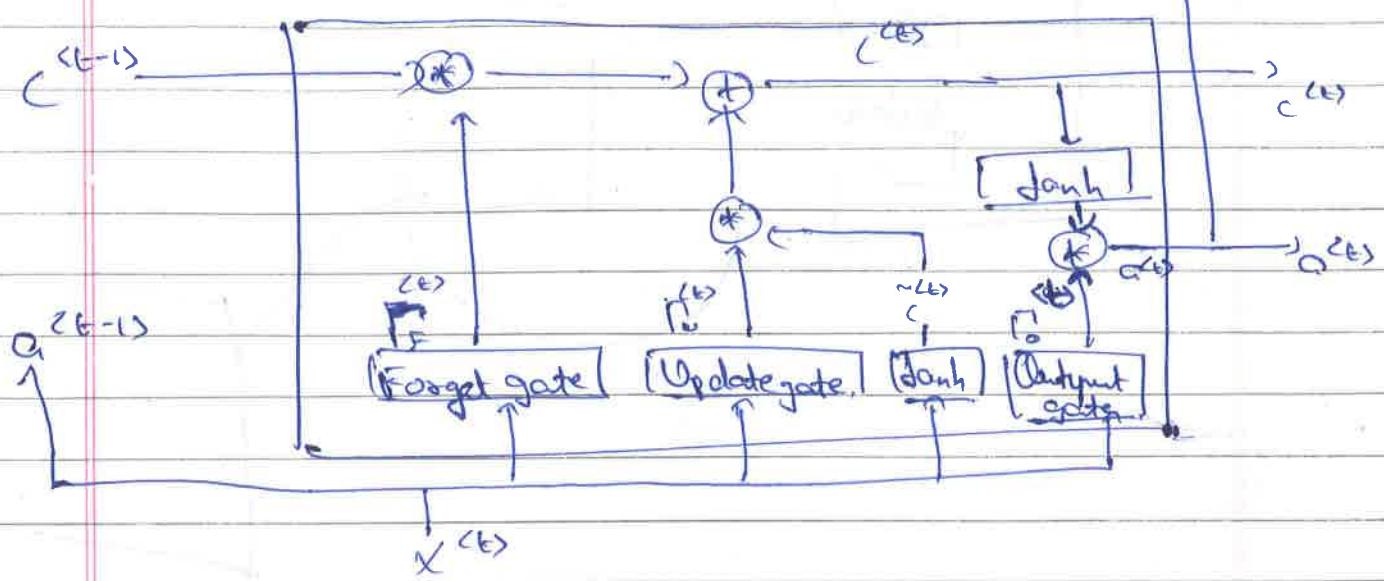
~~forget gate~~  $\Gamma_F = \sigma(\mathbf{w}_f [\mathbf{a}^{(t-1)}, x^{(t)}] + b_f)$

~~output gate~~  $\Gamma_o = \sigma(\mathbf{w}_o [\mathbf{a}^{(t-1)}, x^{(t)}] + b_o)$

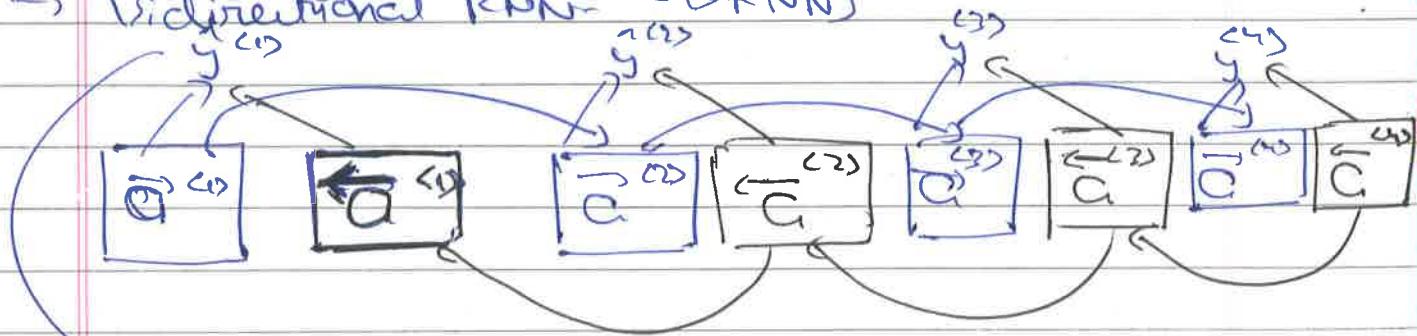
$$c^{(t)} = \Gamma_U * \tilde{c}^{(t)} + \Gamma_F * c^{(t-1)}$$

$$\mathbf{a}^{(t)} = \Gamma_o * \text{Janh} c^{(t)}$$

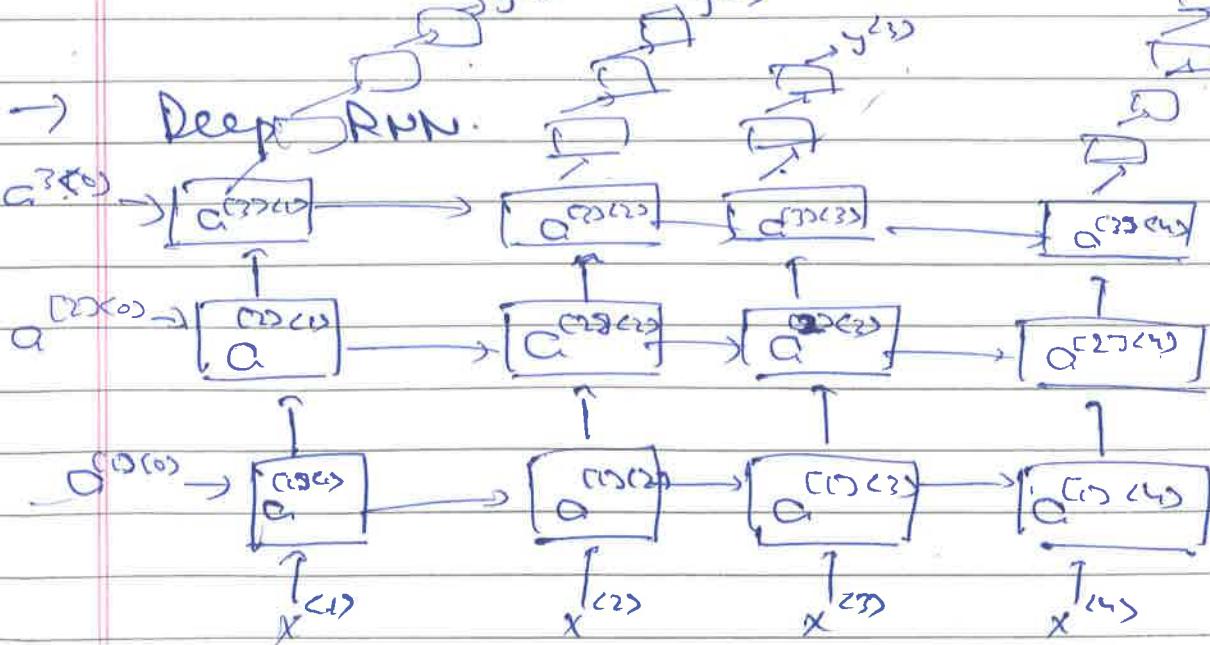
{P-7.0}

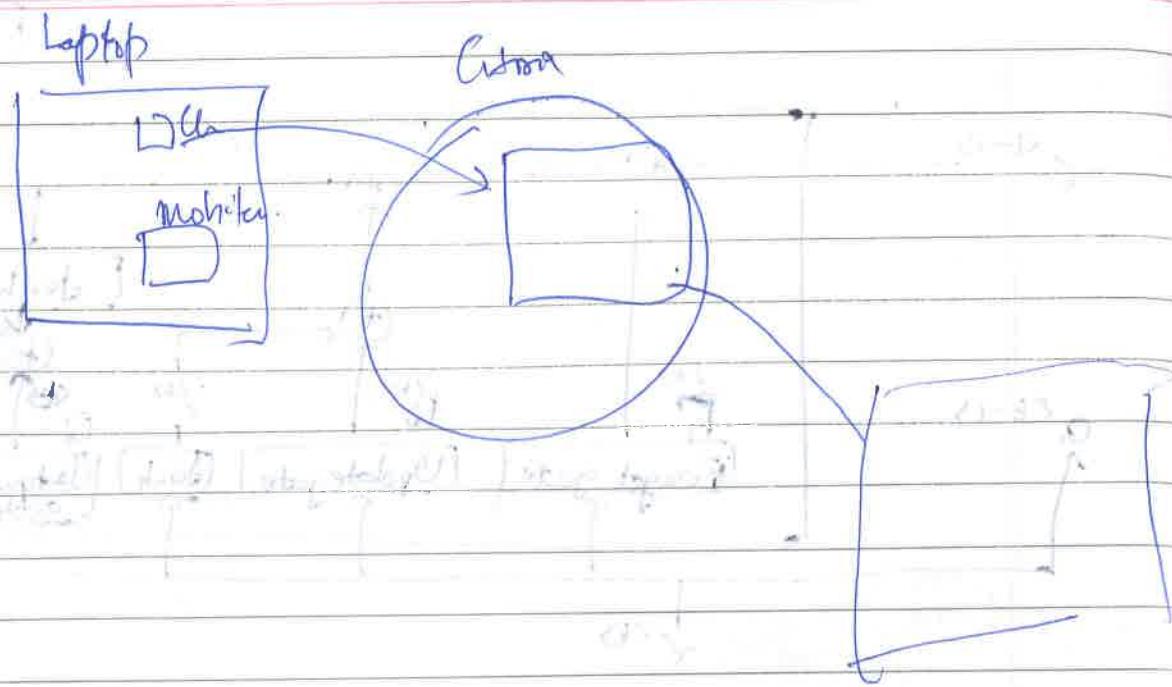


→ Bidirectional RNN  $\Leftarrow$  BRNN



$$\hat{y}^{(t)} = g(w_y[\bar{a}^{(t)}, \bar{c}^{(t)}] + b_y)$$





bash

→ why was nohub  
→ how do you install  
meta

PAGE NO.

DATE

nohub ~~is~~ jupyter notebook &

cat nohub.ctl

ps -ef | grep jupyter

kill -9 982 989 31901

cat nohub.ctl

> nohub.ctl → several entries.

bash script script.sh  
new session → ~~del~~ jps -

du -sh \*

ends at a update

if config

nohub.ctl

on → export  
on → ~~on~~ config  
on → config link  
on → config link

order

022796 44444



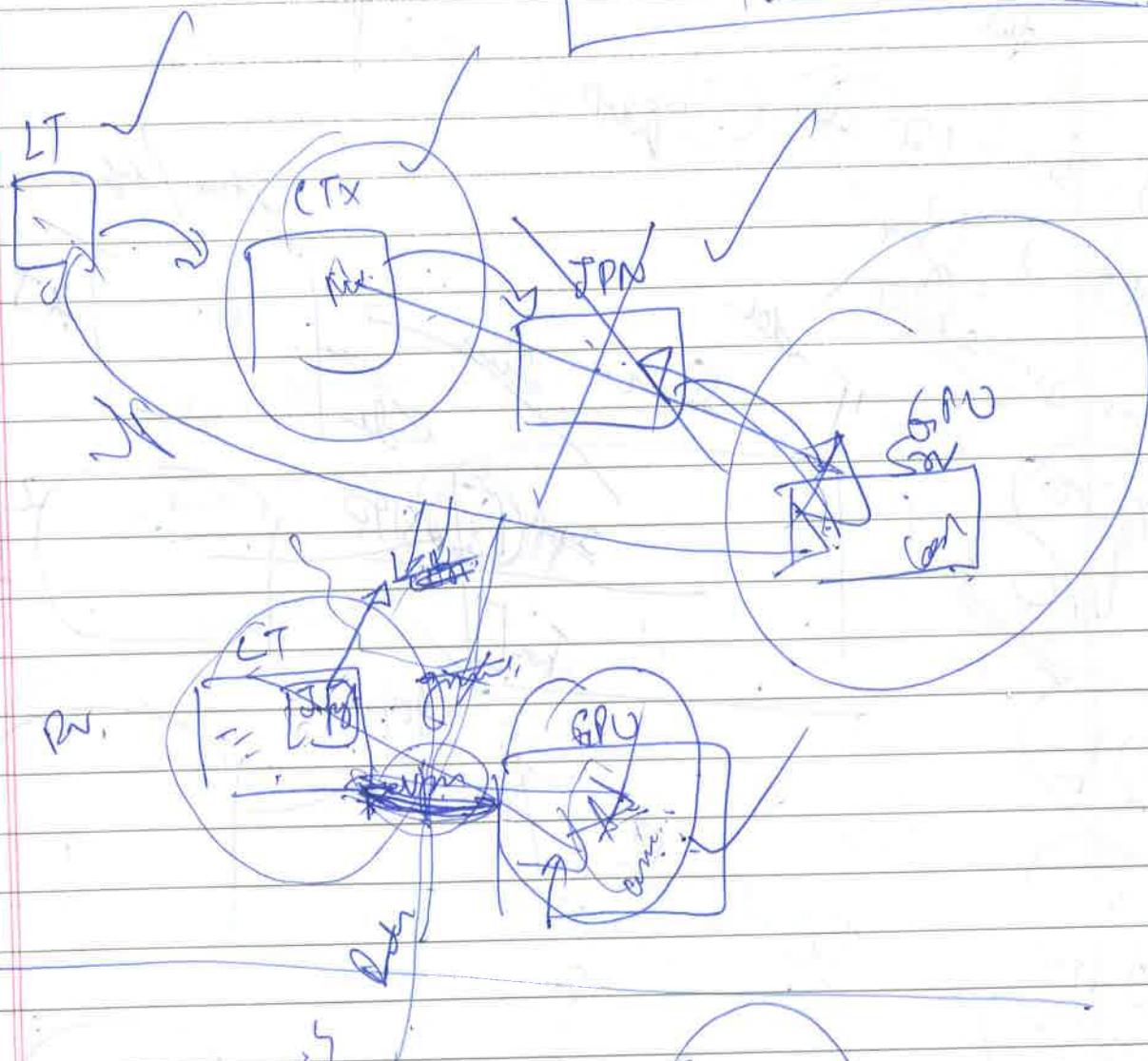
16 146.102.202

100

uhwam:

ifconfig -a

↑ Prog  $m_1 \rightarrow$



b d

$$TD + f(x) = \bar{g}_{(2)}$$

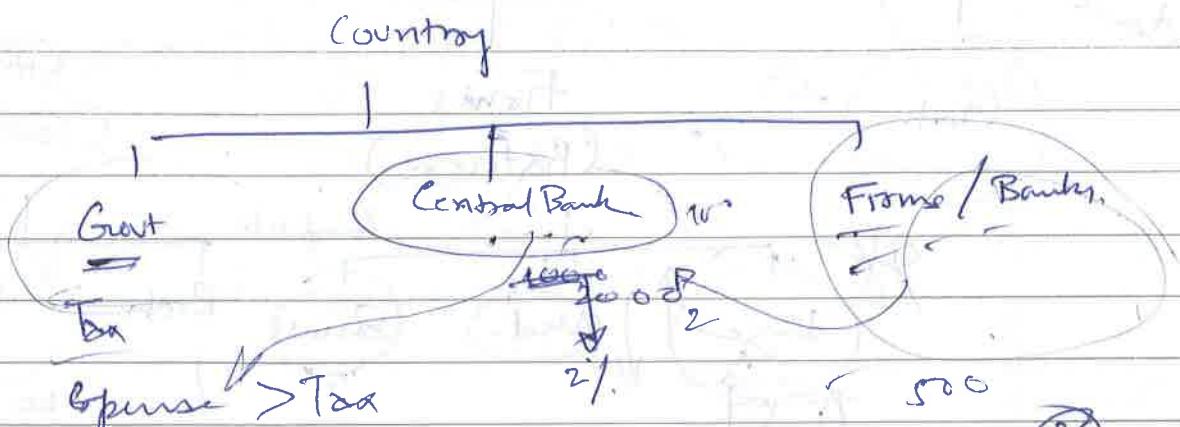
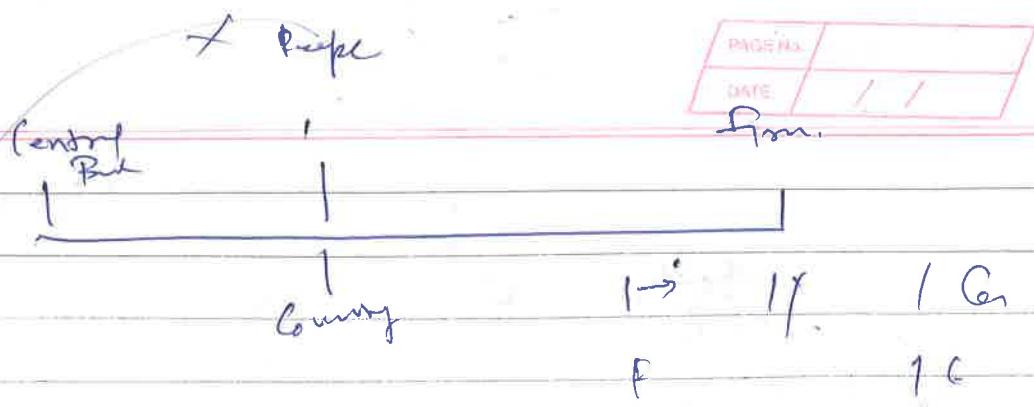
$$P_{\text{dig}} + (f(H) \oplus f(z)) \rightarrow \text{Random generated}$$

~~2008~~ (AP) + 5 (20)

$y = u(x)$

Path

→  $P_{\text{new}}$  → system



1900 → certificate → 5% yield → 100 → 5,  
↓  
15 year

15 200 100

200

2-5 cm.

18 → 40

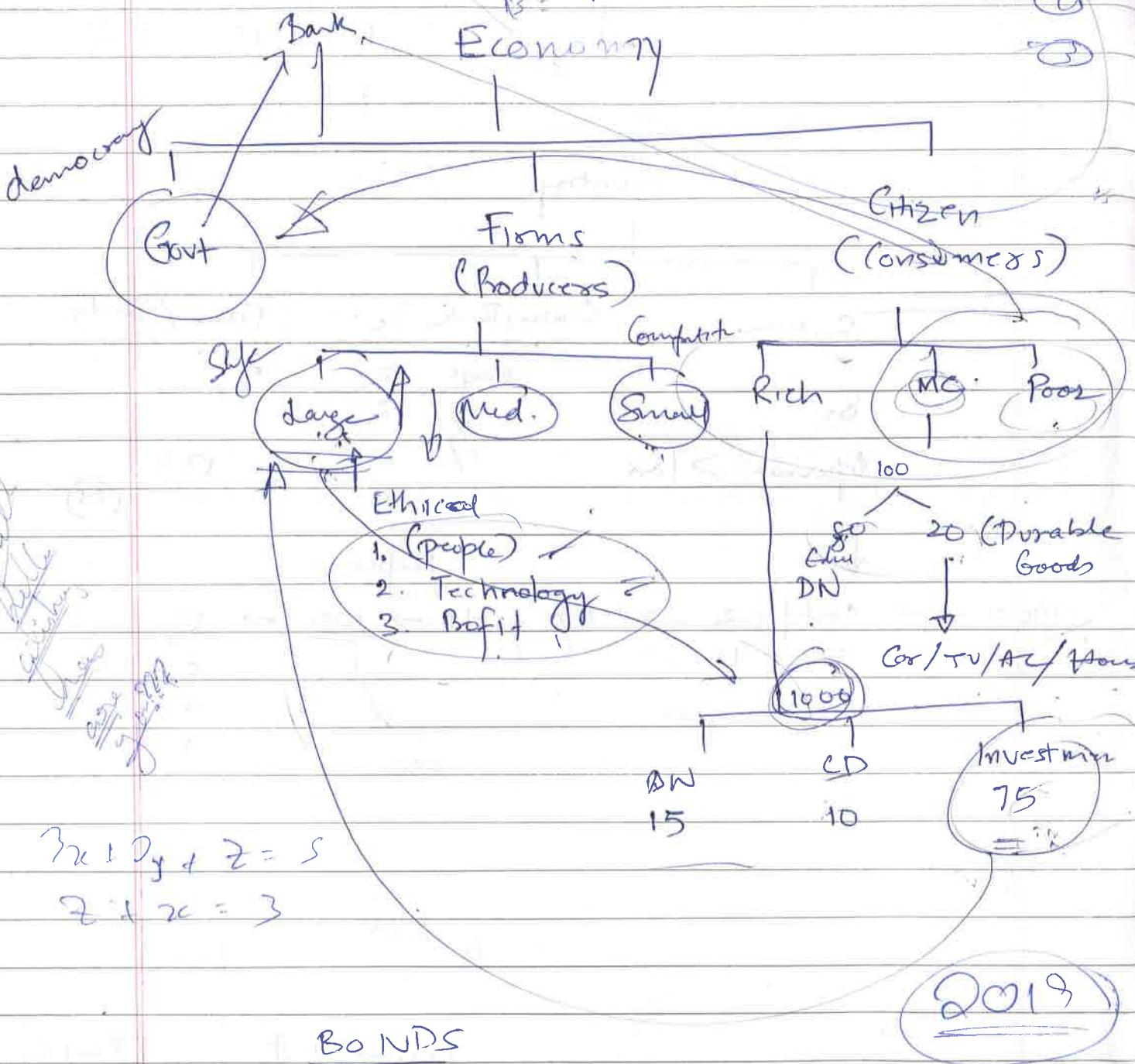
10  
20  
30

20  
10  
0  
0  
0  
0

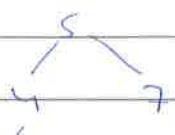
|          |    |
|----------|----|
| PAGE NO. | 11 |
| DATE     |    |

$$x = w_1(x_{w1} + y_{w2}) +$$

$$B =$$



Virtual environment



8P

1937

3

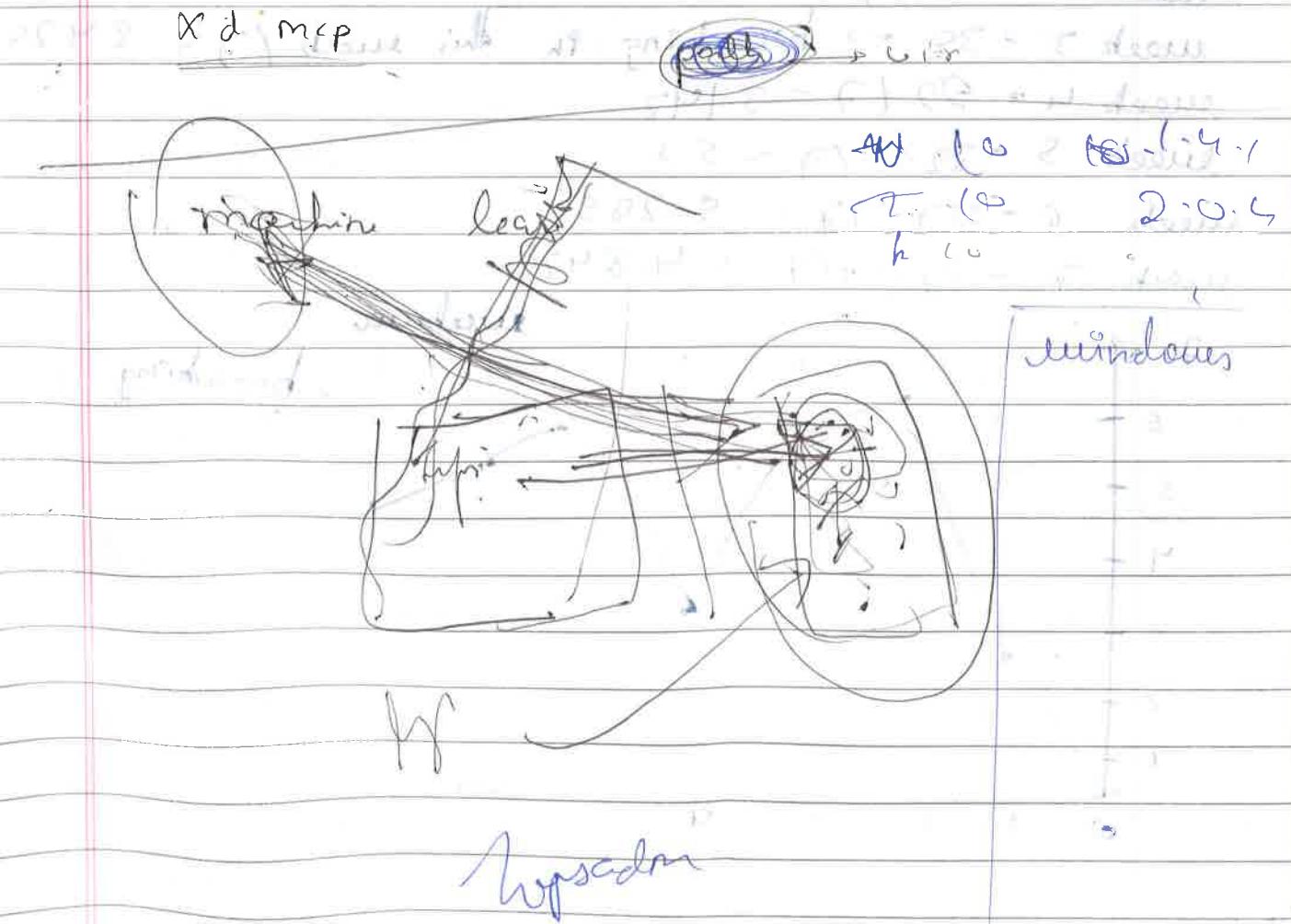
5473

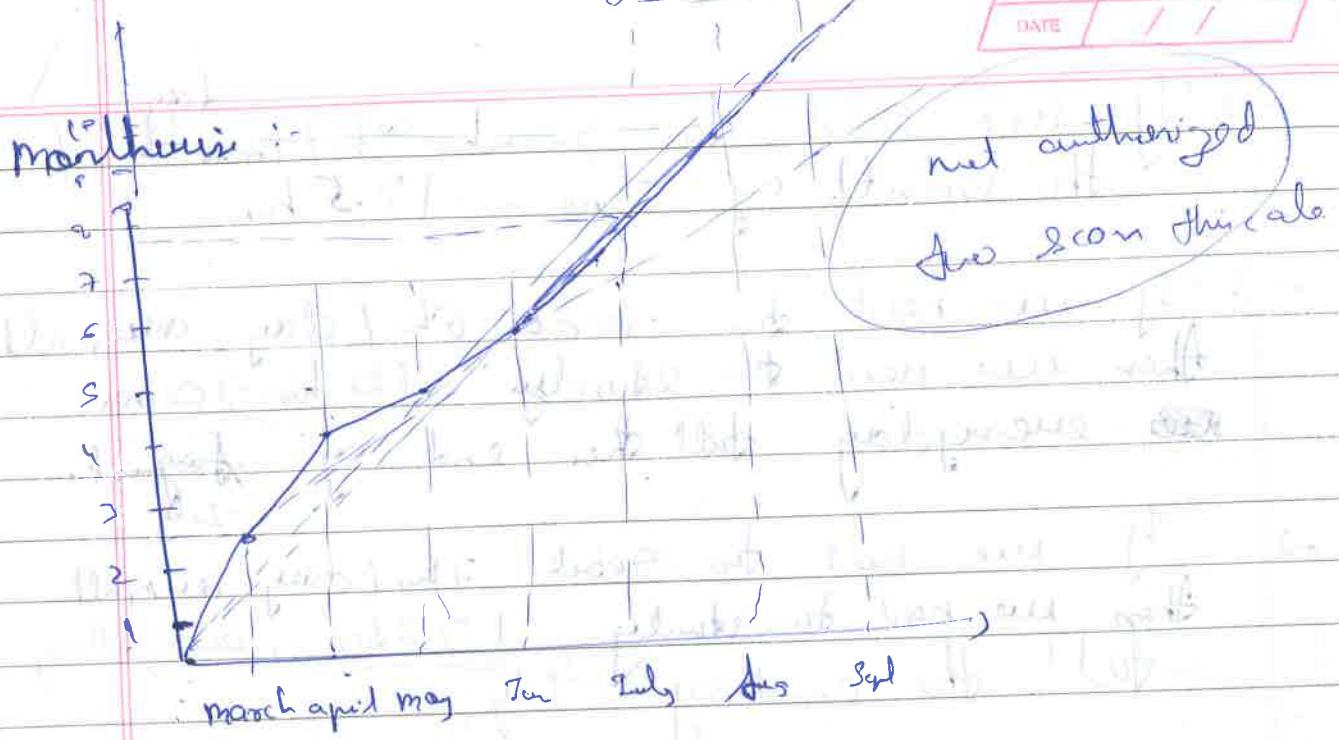
N  
h

|          |  |
|----------|--|
| PAGE NO. |  |
| DATE     |  |

- If we need to reach 9 hr / ~~month~~<sup>day</sup> for the month of June :- 13.5 hr
- If we need to reach 6 hr / day overall then we need to study 10 hr 20 min ~~everyday~~ till the end of ~~August~~<sup>July</sup>.
- If we need to reach 7 hr / day overall then we need to study 17.5 hr everyday till the end of ~~July~~.

Autumn term → Sessions → SEPT → Oct





From may website:-

$$\text{week 1} - 19.5 (7 = 97.85)$$

$$\text{week 2} = 45/7 = 6.4285$$

week 2 - 45 (7) = 0.9200  
week 3 - 38 → breaking in this week (7) = 3.428

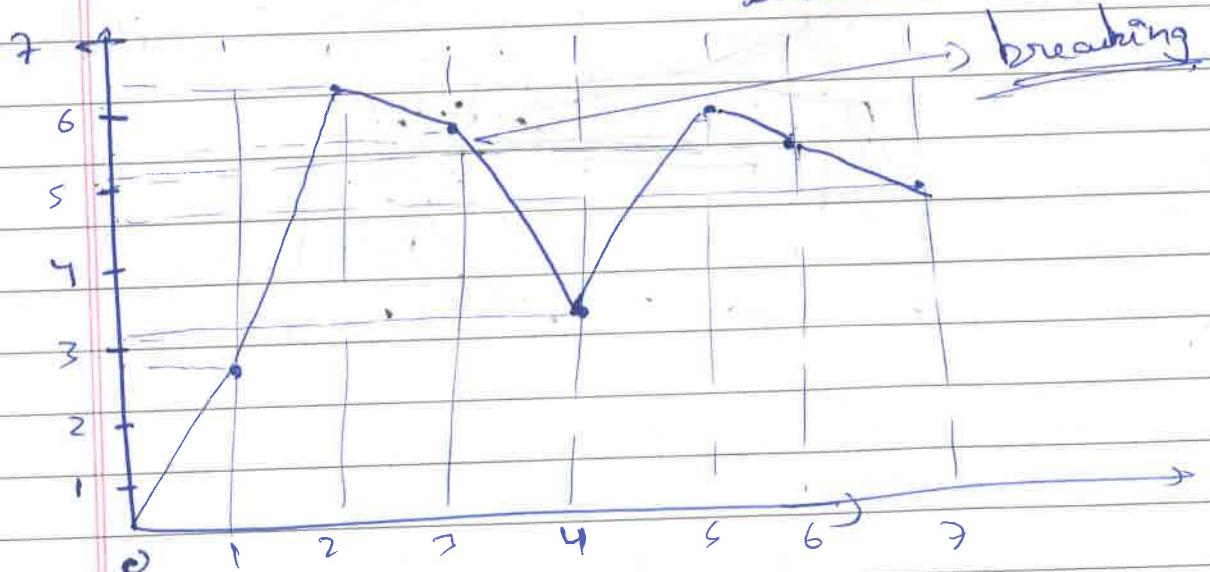
$$\text{week 4} = 2217 = 3.142$$

$$\text{week } S - 38.5 \text{ (7)} = 5.5$$

$$\text{meets } 6 - 3x \text{ } (7) = 5 \cdot 285$$

$$\text{mean } T = 32.5 \text{ (7)} = 4.6428$$

eeeeeeee



finalized with terms, though well, many ~~parameters~~ (76/88)

~~8120~~  $\rightarrow \Delta w[1]$  instead of  ~~$\Delta w[2]$~~  yo 57

(ii) So this is my theory on zero initialization  
of  $w_{s,b}$ . I just wanted to confirm it:-  
→ If we use ~~relu~~ activation function  
[redacted] then the ~~w, b~~ values will be

→ If we use sigmoid activation function then all  
of the nodes of each layer would have the same value  
along with the  $w_{s,b}$  values  
of each layer

and will be updated simultaneously with the same  
values. Thus we could then compress the entire neural  
network to a single equation. (Eg. equation of logistic  
regression)

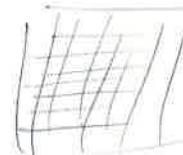
---

discube  $\Rightarrow (A, w, b)$

$\rightarrow A^{[27]}$

cube  $\Rightarrow$  (and linear-cube, activation-cube)

→ We also weight corresponding to node as:  
picture itself with  $((X) \times (Y))$  pixels.



→ In non-conv netwroks the nodes of layers are not necessarily  
direct edge/ corner/ combination of edge. Because even though  
we input a totally random img we still get an output that  
the img is this particular img object with high probability.

normal initialization  $\rightarrow$  rand,