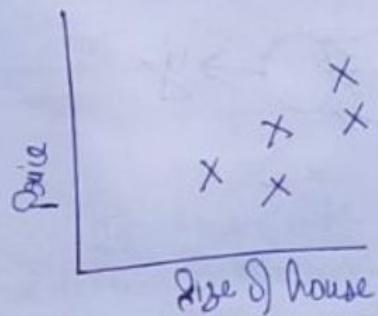


INTRODUCTORY GUIDE TO DEEP LEARNING & NEURAL NETWORK

#1

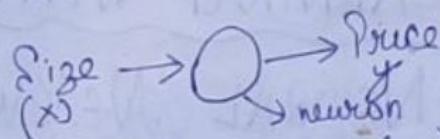
What is Neural Network?

Suppose we want to predict the price of a house. The variables we are given size & house & price.

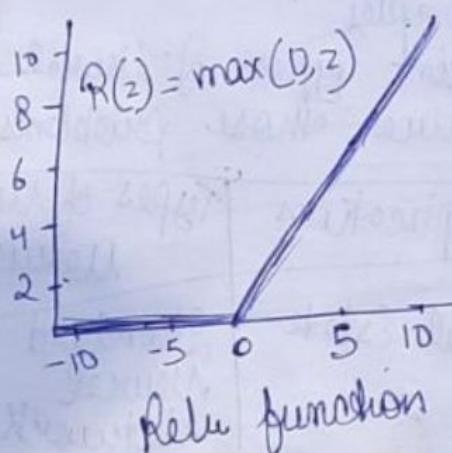


Here, a linear model will try to draw a straight line to fit the data.

So, how to solve this problem using Neural Networks



Here neuron will take an input (size), apply some activation function (ReLU or Sigmoid). One of the commonly used activation function is ReLU (Rectified Linear Unit)

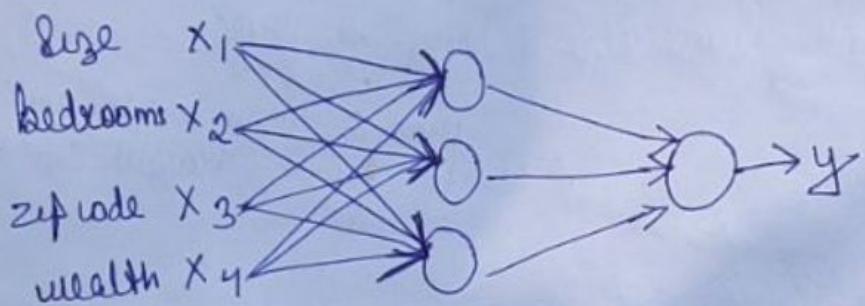


ReLU takes a real number as input and returns the maximum of 0 or that number. So, if we pass 10, the output will be 10, and if the input is -10, the output will be 0.

We'll consider multiple features like number of bedroom, size, postal code.

Here, we have pass 4 features as input to the neural network.

This is how a neural network with 4 inputs and an output with single hidden layer will look like:



SUPERVISED

you have input x and you want to learn a function
→ mapping to some output.

LEARNING WITH NEURAL NETWORKS

Supervised learning refers to find a function that can map input and its corresponding output. We have to define output for each given input and we train the model on these examples.

Following are the different types of Neural Network that can be used to solve those problems.

Input	Output	Application	Types of Neural Network
Home features	Price	Real Estate	Standard Neural Network

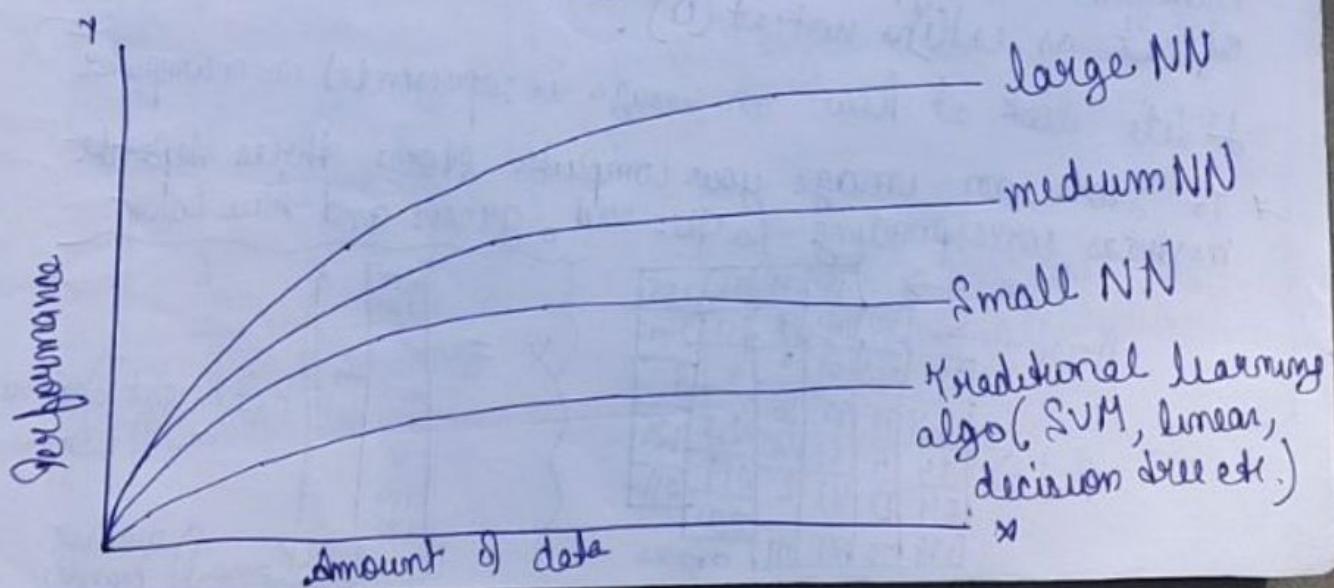
②	Ad, user info	Click prediction (1/0)	Online Advertising	Standard Neural Network
③	Image	Image class	Photo Tagging	CNN
④	Audio	Image class	Photo Tagging Speech Recognition	RNN
⑤	English	Chinese	Machine Translation	RNN
⑥	Image, Radar info Information	Position of car	Autonomous Driving	Custom/ Hybrid NN

Supervised Learning

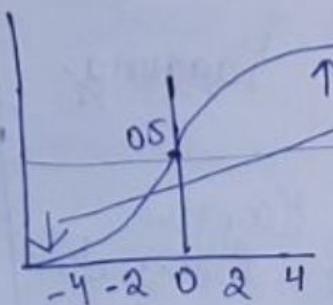
structured data
⇒ data in table form

unstructured data
⇒ audio
⇒ Text
⇒ Video
⇒ Picture

WHY DEEP LEARNING PROGRESS TAKING OFF?



To improve the computation time of the model, activation function plays an important role. If we use sigmoid function



Where we can see here with the slope of the function the gradient is nearly zero and so learning become very slow.

The slope or gradient of this function, at the extreme ends is close to zero. Therefore, the parameters are updated very slowly, resulting in very slow learning. Hence switching from sigmoid to ReLU (Rectified Linear Unit) is one of the biggest breakthroughs we have seen in Neural Network. This is the primary reason for faster computation of the model.

LOGISTIC REGRESSION AS A NEURAL NETWORK

① BINARY CLASSIFICATION

BINARY CLASSIFICATION

Logistic Regression is an algorithm for binary classification. Suppose we have to recognise a photo of cat as ~~cat~~ not cat (0).

A person's image is represented in a wmp

of cat as $\text{cat}(\theta)$ is not cat(0).
Hence look at how an image is represented in a computer
image your computer stores three separate

- * To store an image your computer stores three separate matrices corresponding to the red, green and blue color.

	Blue	Green	Red	Blue	Green	Red
	→	→	→	150	175	11
				200	200	7
				255	255	9
				255	255	8
				123	97	19
				95	88	219
				88	80	200
				54	41	111
				119	13	111

So if your input image is 64 pixels by 64 pixels then we'll have

$3 \times 64 \times 64$ matrices.

↳ Red, Green, Blue.

What we are going to do is define a feature vector x corresponding to this image. see picture

NOTATION

A single training example is represented by a pair of (x, y) where x is an n_x -dimensional feature and y the label is either 0 or 1 (according to binary classification).

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

$n_x = \text{no. of rows}$
Here we have only
a column 64×1

$m = \text{no. of training examples}$

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

and your training sets will be written (x_1, y_1) which is the input and output for your first training example $(x^{(1)}, y^{(1)})$ from above up which is upto $(x^{(m)}, y^{(m)})$

↳ all together is your training set

$M_{train} = \text{No. of training examples}$
 $M_{test} = \text{No. of test examples}$

We are going to define a matrix, capital X . As defined by taking your training set input x_1, x_2 .

$x_1 = \text{first column}$
 $x_2 = \text{second column}$

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

$n_x = \text{height of matrix}$
 $= \text{no. of rows}$

$m = \text{no. of training examples}$

We can also use transpose of $X (X^T)$ but it will make our calculations more complex.

$$X = \mathbb{R}^{n \times m}$$

n : number of rows
 m : number of columns

Now what about the output labels Y ?

To make your implementation of a neural network easier it would be convenient to also stack Y in columns

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

$$Y \in \mathbb{R}^{1 \times m} \xrightarrow{\text{no. of columns}} \boxed{Y. \text{shape} = (1, m)}$$

\xrightarrow{\text{no. of rows}}

② Logistic Regression

formula for logistic regression is

$$\varphi(x) = \frac{1}{1 + e^{-(ax + b)}} \quad a = \text{slope}, \ b = \text{intercept}$$

$e = \text{euler's}$

x = matrix of df

y = target out of datframe

\hat{y} = (y hat) = predicted output

w = weight

b = bias

what are weights?

different weights are given to neurons associated with input layer. Random weights are given in forward propagation and weights changed during backward propagation.

what are bias?

Bias are like intercept in linear regression they are the values which are added in activation function. During initializing process random weights and bias are put up.

To Output
To predictive output for sigmoid function @ is

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

$z = w^T x + b$

↓ sigmoid function

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

b = bias
 x = Input metrics
 w^T = transpose of weight.

③ LOGISTIC REGRESSION COST FUNCTION

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

here λ = loss

\hat{y} = predicted y value

y = Actual output

To obtain the parameters w &
B of the logistic regression
model, we need cost function

loss function is a function measure how good our output (\hat{y})
is when the true label y (actual y value)
~~output approx~~

we can also define our loss function as $\frac{1}{2}(\hat{y} - y)^2$
but this will give us an optimization problem that is convex.

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

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

Here, ^{This is only for binary classification. We use cross-entropy for multiclassification. Cross entropy → multiclassification}

If $y=1$ then we want \hat{y} as large as possible

If $y=0$ then we want \hat{y} as small as possible

Note: the loss function ^{was} defined here is with respect to a single training example. It measures how well you're doing on a single training example

Cost function

* It measures how you're doing on entire training set. So cost function J applied to your parameter W (weight) and B (bias).

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

(Applied to each i) training example

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

* loss function is applied to just single training example & the cost function is other cost of the parameter. So in training your logistic regression model we're going to try to find parameters w & b that minimizes the overall cost function.

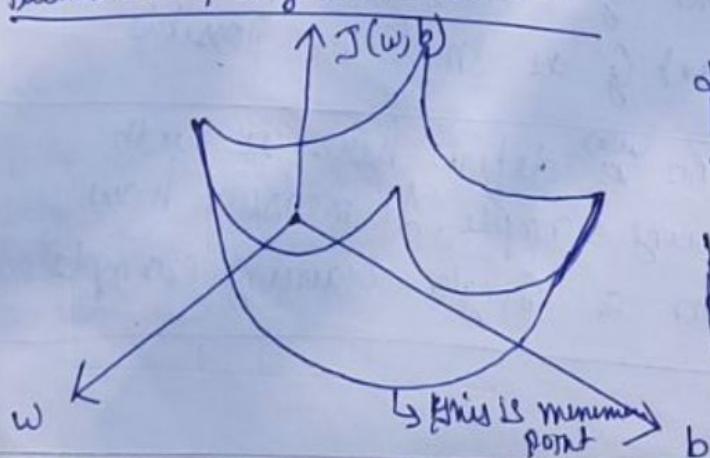
④ GRADIENT DESCENT

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

Here it's natural that we want to minimize the cost function.

Illustration of gradient descent

(This is convex function)



the height of this graph is $J(w, b)$

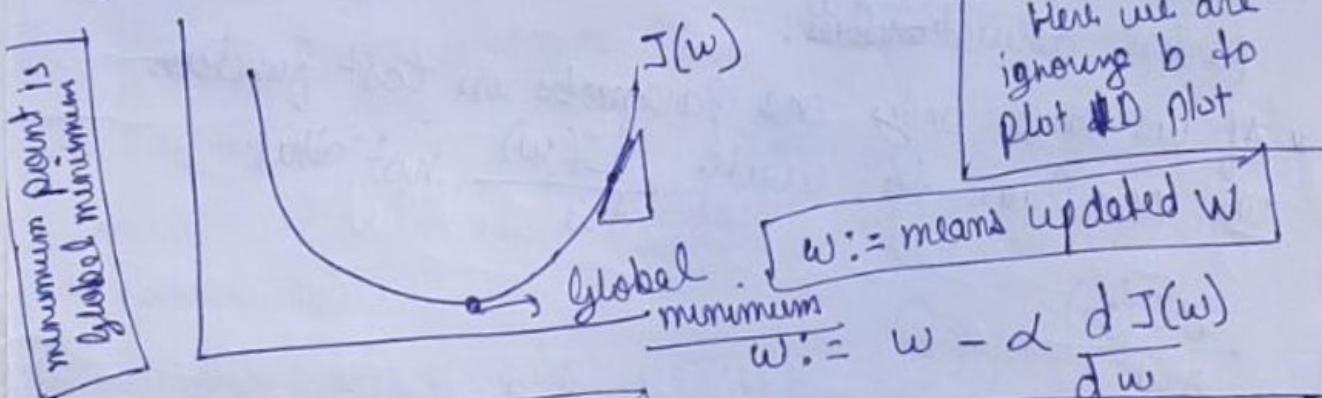
Here we want to find the value of w & b which is minimum

Note: derivative = slope = gradient descent

- So, to find a good value for the parameter what we'll do is initialize w and b to some initial value. For logistic regression any initial value will work.
- * and what gradient descent does is it starts at that initial point and then takes a step in the steepest downhill direction.
 - * So after taking some steps you'll reach the steepest point. (at minimum value) \hat{w} and \hat{b}

let's go into bit details

- * For the purpose of illustration, let's say that there's some function $J(w)$ that you want to minimize



So, the slope of the function is really the height divided by the width.

Now, we wrote our gradient descent for $J(s)$ if only w was your parameter

Note: $\frac{d}{dw}$ are called derivative

$\alpha = \text{alpha} = \text{learning rate}$
 $\alpha =$ It tells how big a step we take on each gradient descent

$\frac{d}{dw}$: This is basically the update or the change we want to make to parameters w .

Note: when we'll write the code you write " dw " in the code which represents $\frac{d}{dw}$.

In logistic regression, your cost function is a function of both w and b ($J(w, b)$)

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

we drop

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

~~b is fixed~~

$w :=$ means updated weights

$b :=$ means updated bias

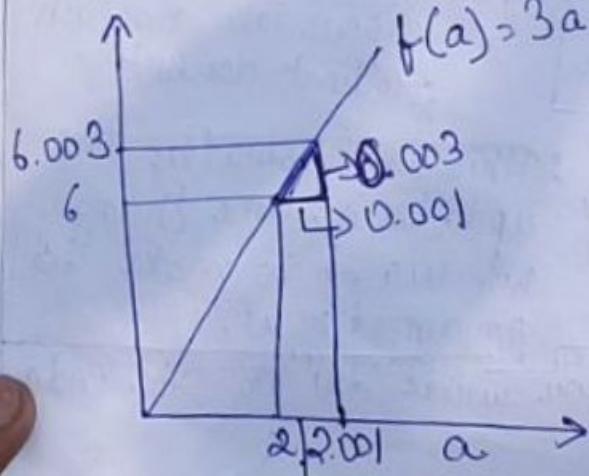
Some derivatives information

- * When we have α parameters in our cost function we have to write $\frac{\partial J(w, b)}{\partial w}$ instead of this $\frac{\partial J(w, b)}{\partial w}$. This is called partial derivatives.

- * If we have only one parameter in cost function we can have to write $\frac{\partial J(w)}{\partial w}$ not this

$$\frac{\partial J(w)}{\partial w}$$

⑤ DERIVATIVES



If we put value of $a = 2$
we'll get $f(a) = 6$
at $a = 2$, $f(a) = 6$

at $a = 2.001$, $f(a) = 6.003$
As we see
height divided by width

Height divided by width

Slope (derivative) of $f(a)$ = $\frac{\text{height}}{\text{width}} = \frac{\text{change in } y}{\text{change in } x}$

\Rightarrow This means that when we add 0.001 to the x axis y will increase by 0.003.
With video for complete

$$= \frac{0.003}{0.001} = 3 \quad (\text{here slope is 3})$$

So this is what this equation means ($f(a) = 3a$)

at $a = 5$

$a = 5.001$

$f(a) = 15$

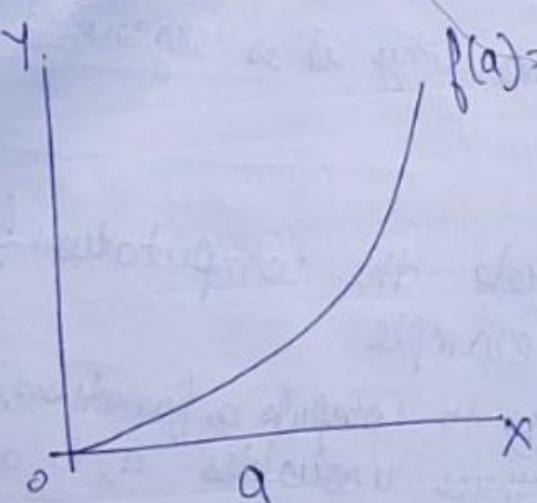
$f(a) = 5.003$

So, the way we wrote this, that the slope of the function f is equal to three is

$\frac{df(a)}{da}$ what this means is the slope of a function $f(a)$ when we add a with a tiny little amount
 ↓
 tiny little amount
 this is equal to 3

we can also write $\frac{df(a)}{da}$ as $\frac{d}{da} f(a)$

⑥ MORE DERIVATIVE EXAMPLES



$f(a) = a^2$

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

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



How?

because here line is not straight because of this values of after decimals changes
 Here is how

$$\begin{array}{l} \cancel{a=2.00} \\ \cancel{a=2.001} \end{array}$$

$$\begin{array}{l} f(a)=4.004 \\ f(a)=a^2 \end{array}$$

~~a=2.001~~

$$\begin{array}{lll} \textcircled{1} \quad f(a) = a^2 & \frac{d}{da} f(a) = 2a & \text{for } a = 2.001 \\ & & f(a) = 2a \\ & & \downarrow \\ & & f(a) = 4.004 \end{array}$$

lets take another example to understand this clearly

$$\begin{array}{lll} \textcircled{2} \quad f(a) = a^3 & \frac{d}{da} f(a) = 3a^2 & a = 2.001 \\ & 6a & f(a) = 8.012 \end{array}$$

⑦ COMPUTATION GRAPH

Neural networks are organised in terms of a forward pass (forward propagation) followed by backward pass (back propagation)

Computation graph explains why it is organized this way.

Computations

In order to ~~comp~~ illustrate the computation graph let's take a simple example

lets say that we're trying to compute a function, J , which is the function of three variables a , b and c

lets say

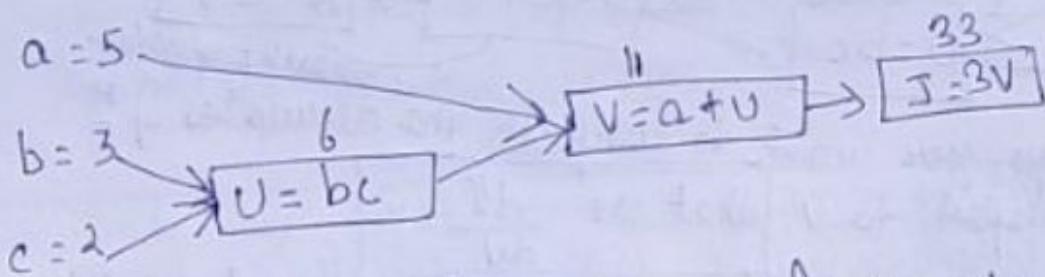
$$J(a, b, c) = 3(a + bc)$$

$$U = bc$$

$$V = a + U$$

$$J = 3V$$

lets take a example here



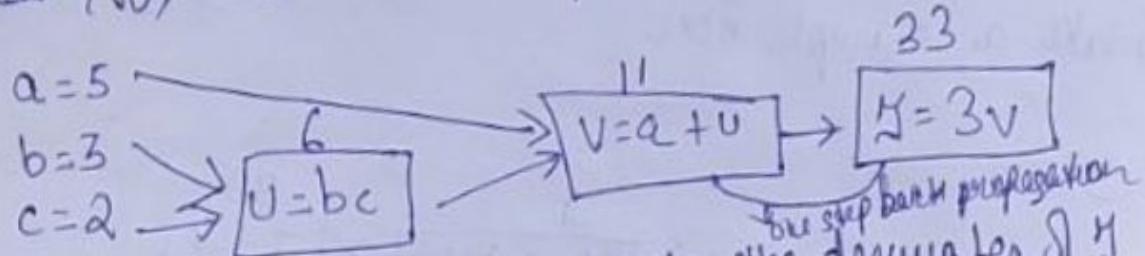
So, the computation graph comes in handy when there is some distinguished or some special output variable such as J in the case that you want to optimize.

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

Computation Graph

⑧ DERIVATIVES WITH COMPUTATION GRAPH

So, here's a computation graph



Let's say you want to compute the derivatives of Y with respect to V . That is $\frac{dY}{dV} = ?$

What this means that if we take this value δY and change it a little bit how would the value of Y change.

$$Y = 3V \quad \text{according to formula}$$

$$\text{and right now } V = 11$$

So if we want to increase a little bit value of V from 11 to 11.001 then Y will jump up from 33 to 33.003.

So the derivative of Y with respect to V is equal to 3
that is $\frac{dY}{dV} = 3$ because the increase in Y is 3 times the increase in V . $\frac{dY}{dV} = 3 \frac{dV}{dV} = 3$

In one terminology of backpropagation is that if you want to compute the derivatives of this final output variable with respect to V , then we've done one step of backpropagation

Now, what's the value of $\frac{dY}{da} = ?$

Note: when we calculate $\frac{dY}{da}$ then this is one step of propagation

$\frac{dJ}{da} = ?$ $a = 5$ (for our example) [If we pump up the value of a , how does that effect the value of J ?]

$v = a + u = 11$ (for our example)

Not diff a

- If we increase the value of a by from 5 to 5.001
 v which was $a+u$. This would get increased to 11.001
and also J moves up from 33 to 33.003
This means this derivative ($\frac{dJ}{da} = 3$)

Note: One way to break this down is to say that if you change a then it will change v and through changing v that would change J .
In calculus this is called Chain rule.

- If you increase a by 0.001 v changes by the same amount. So, $\frac{dv}{da} = 1$

Wrapping up:-

Here: $\frac{dv}{dJ} = 3$ and $\frac{dJ}{da} = 1$

So the product of these is 3 times 1.

So, that's gives you a the correct value of $\frac{dJ}{da}$
 $= 3$.

$$\frac{dJ}{du} = \frac{dJ}{dv} \cdot \frac{dv}{du}$$

multiplication

$$3 = 3 \times 1$$

$$3 = 3$$

$$d \frac{dy}{db} = \frac{dy}{du} \cdot \frac{du}{db}$$

$$= 3 \times 2$$

$$db = 6$$

$c=2$

d

$$* \frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc}$$

$$= 3 \cdot b$$

$$= 3 \times 3$$

$b=3$

$dc = 9$

④ LOGISTIC REGRESSION GRADIENT DESCENT

Logistic regression help

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

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

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

Logistic regression derivatives

→ Here we have only 2 feature x_1, x_2

$\begin{array}{c} \text{z} = w_1x_1 + w_2x_2 + b \\ \text{a} = \sigma(z) \\ \text{d}z = \frac{d\text{z}}{dz} = \frac{d\text{d}(a,y)}{da} \end{array}$ → explained below

$\begin{array}{c} \text{d}z = \frac{d\text{z}}{dz} = \frac{d\text{d}(a,y)}{dz} \\ = a - y \\ = \frac{da}{da} \cdot \frac{da}{dy} \end{array}$ → $\begin{array}{c} \text{d}a \quad [\text{python}] \\ \text{d}a = \frac{da}{dy} \\ = \frac{y}{a} + \frac{1-y}{1-a} \end{array}$

(So in logistic regression we want to modify our parameters)

$$\begin{aligned} \frac{\partial \text{J}}{\partial w_1} &= "dw_1" = x_1 \cdot d_z \\ dw_2 &= x_2 \cdot d_z \\ db &= d_z \end{aligned}$$

GRADIENT DESCENT ON m Examples

GRADIENT DESCENT

LOOPS OVER ALL TRAINING EXAMPLES

In the previous topic, you saw how to compute derivative and implement of gradient descent with respect to just one training example. In this we just implement in all training examples.

we can define the prediction and cost function for 'm' training examples as:

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

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

~~The derivative of a loss function have to be apply Logistic regression for this example~~

The derivative of a loss function can be written as :

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

From the above we can conclude the logistic regression code:

$$z = 0; dw_1 = 0; dw_2 = 0, db = 0; \# \text{Dev}$$

$$w_1 = 0; w = \theta, b = 0;$$

② for $i = 1$ to m

For-ward pass

$$z(i) = w_1 * x_1(i) + w_2 * x_2(i) + b$$

$$a(i) = \text{Sigmoid}(\sigma(z(i)))$$

$$\Delta J = (y(i) * \log(a(i)) + (1 - y(i)) * \log(1 - a(i)))$$

Backward pass

$$dz(i) = a(i) - y(i)$$

$$dw_1 += dz(i) * x_1(i)$$

$$dw_2 += dz(i) * x_2(i)$$

$$db += dz(i)$$

$$J / = m$$

$$dw_1 / = m$$

$$dw_2 / = m$$

$$db / = m$$

Gradient descent

$$w_1 = w_1 - \alpha \gamma * dw_1$$

$$w_2 = w_2 - \alpha \beta * dw_2$$

$$b = b - \alpha \delta * db$$

PYTHON AND VECTORIZATION

① Vectorization

vectorization is basically the art of getting rid of explicit loops in your code. It executes the code fast.

$$z = w^t x + b \quad [w \text{ and } x \text{ are vectors}]$$

Non-vectorization takes more time than vectorization

Non-vectorized implementation Vectorized implementation

$$z = 0$$

for i in range(n, x)

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

$$z += b$$

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

Code implementation for both together so that we can see some taken by each of them when we run

```
In [1]: import numpy as np
       a = np.array([1, 2, 3, 4])
       print(a)
```

```
Out[1]: [1, 2, 3, 4]
```

import dune

a = np.random.rand(1000000) → It creates an array of specified shape and fills it with random values
b = np.random.rand(1000000)

tic = time.time() → some taken to implement this code

c = np.dot(a, b) → $\underline{ax + b}$

toc = time.time()

print(t)

print("Vectorized version:" + str(1000 * (toc - tic)) + "ms")
↳ here we have multiplied by 1,000 to get milliseconds

c = 0

tic = dune.dume()

for i in range(1000000):

c += a[i] * b[i]

toc = time.time()

print(t)

print("For Loop:" + str(1000 * (toc - tic)) + "ms")

Out

250286.989866

vectorized version: 1.5027523... ms

250286.989866

for loop: 474.29513931.. ms

So, the non-vectorize version took something like 300 times longer than the vectorize version.

② MORE VECTORIZATION EXAMPLES

[whenever possible, avoid explicit for-loops]

If you want to compute a vector U as the product of matrix A and another vector V

$$U = AV$$

But this is define as $U_i = \sum A_{ij} V_j$

Implementation of this in numpy will be

$$u = np.\text{dot}(A, v) \rightarrow \underline{\text{this is vectorized}}$$

Vectors and matrix valued functions

Say you need to apply the exponential operation on every element of a matrix/vector:

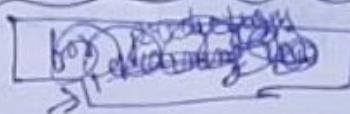
Exponential function have the form $f(x) = b^x$, where $b > 0$ and $b \neq 1$

A. `import numpy as np
u = np.exp(v)`

↳ If we apply this as non-vectorized in python i.e. have use of for loops and for loops are time consuming

* Whenever you are tempted to write a for loop take a look, and see if there's a way to call to call a Numpy built-in function to do it without the for-loop

③ VECTORIZING LOGISTIC REGRESSION



LOGISTIC REGRESSION

$$\begin{aligned} z^{(1)} &= w^t x^{(1)} + b & z^{(2)} &= w^t x^{(2)} + b & z^{(3)} &= w^t x^{(3)} + b \\ a^{(1)} &= \sigma(z^{(1)}) & a^{(2)} &= \sigma(z^{(2)}) & a^{(3)} &= \sigma(z^{(3)}) \end{aligned}$$

we might need to do this m times if we have m training examples

so how we can calculate all the value of $z^{(1)}$ without using for loop in python.

let's see how we can do it.

$$X = \underset{\text{Training input}}{\downarrow} \quad W^T W b \quad \left[\begin{array}{c} X^{(1)} \\ X^{(2)} \\ \vdots \\ X^{(m)} \end{array} \right] \quad \begin{array}{l} \text{thus is} \\ (h \times m) \text{ matrix} \end{array}$$

How do calculate $z^{(1)}$ all at same time?

$$\left[z^{(1)} \ z^{(2)} \ \dots \ z^{(m)} \right] = \underset{\substack{\text{Training inpt} \\ \downarrow}}{W^T X} + \left[\underset{\substack{\text{bias} \\ \downarrow}}{b \ b \ b \ b \ b} \right]$$

W^T will be a new vector

This is equal to equation on starting of topic

An numpy we can calculate

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

as

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

we'll end up
adding b at every
element

\hookrightarrow b is 1×1 matrix [just
one @ number]

④ BROADCASTING IN PYTHON

Broadcasting is another technique that you can use to make your Python code run faster.

Example:

Calories from carb, protein, fat in 100g of different foods:

	Apples	Beef	Eggs	Potatoes
Carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	135.0	99.0	0.9

But what if I want to calculate Average of these numbers in python

In | 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, 135.0, 99.0, 0.9]])

print(A)

Out

In | A.sum (axis=0)

print (A.sum)

Out | [59. 239. 155.4 76.9]

In | percentage = 100 * A / A.sum (axis=0)

print (percentage)

we don't have to use reshape here because our matrix (el) is already (1,4)

Out

| [[- , - , - , -]]
| [- , - , - , -]]
| [- , - , - , -]]

Broadcasting more examples

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

don't use this

$$\text{Q2} \quad \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ (m, n) \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 203 & 303 \\ 104 & 205 & 306 \\ 100 & 200 & 300 \end{bmatrix}$$

⑤ A NOTE ON NUMPY VECTORS

~~TOP~~ Tips and tricks to eliminate all the strange looking bugs in code :-

* when you want to create a random matrix of 5 row and 1 column we generally write
 $a = np.random.randn(5)$

but when we print this variable a will get 5 column instead of 5 row.

So always use this

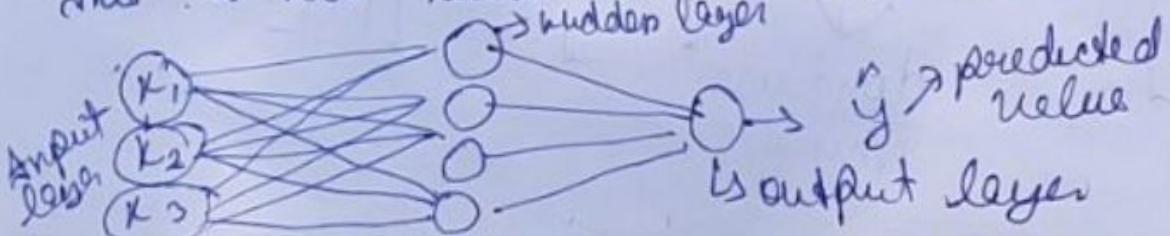
$a = np.random.randn(5, 1)$

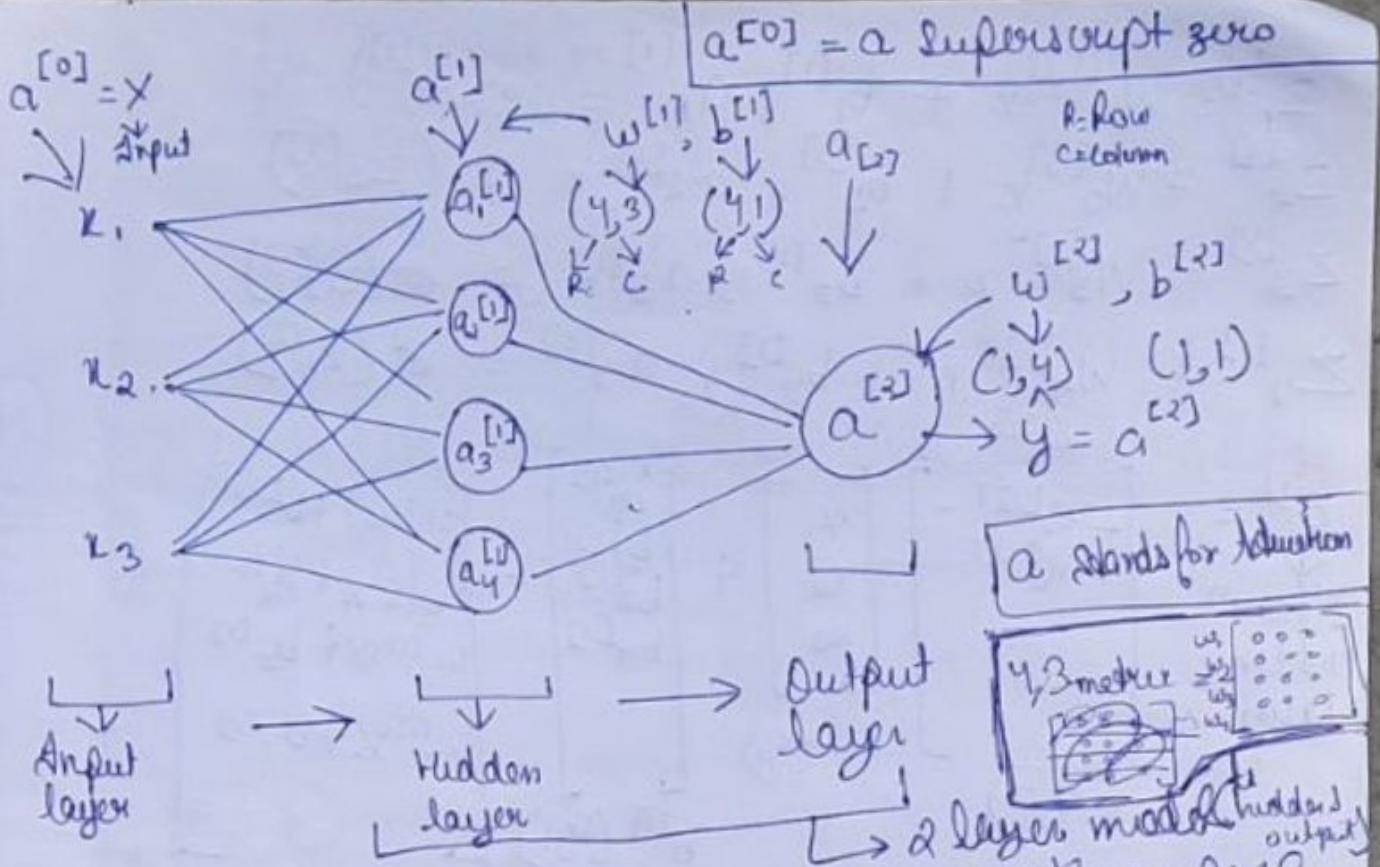
Note:- Don't use rank 1 array

WEEK 3

① NEURAL NETWORK REPRESENTATION

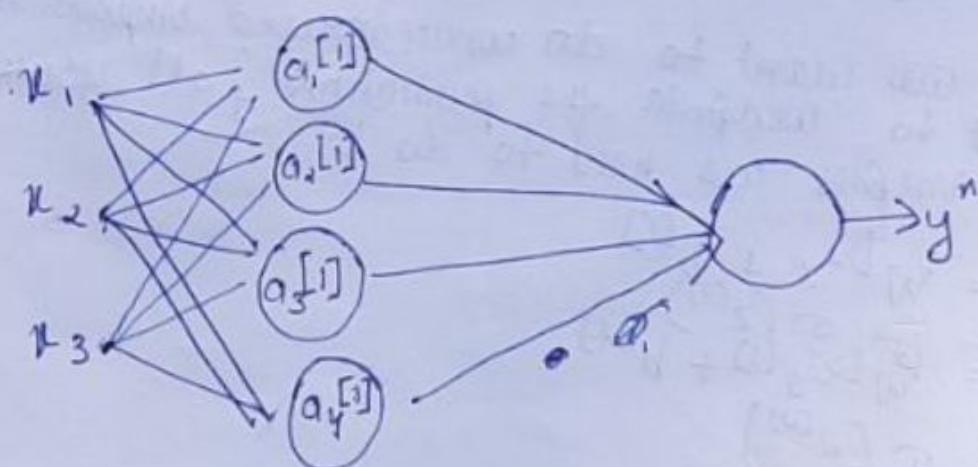
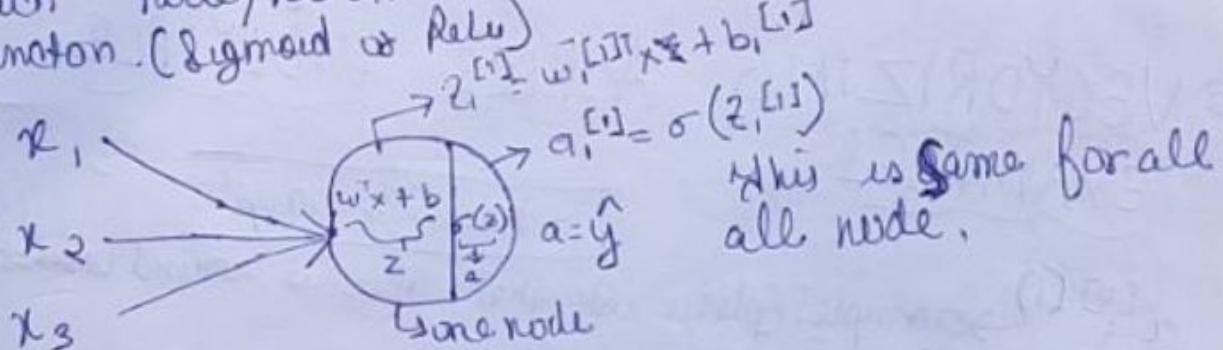
This is how neural network looks like:-





② COMPUTING A NEURAL NETWORKS OUTPUT

Each node/neuron in the hidden has Activation function. (Sigmoid or ReLU)



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

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

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

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

\downarrow
 z^[1] or 1st
 Hidden
 layer

$$\begin{aligned}
 z^{[1]} &= \begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \\ \vdots \\ -w_4^{[1]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ w_3^{[1]T} x + b_3^{[1]} \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix} \\
 &\quad = \begin{bmatrix} t_1^{[0]} \\ t_2^{[0]} \\ t_3^{[0]} \\ t_4^{[0]} \end{bmatrix}
 \end{aligned}$$

③ VECTORIZING ACROSS MULTIPLE EXAMPLES

$a^{[0]}(i)$ → example_i (please note that this is round bracket)
 → layer 2

i = training example

Suppose we want to do unvectorized implementation
 and want to compute the predictions of all your
 training examples we need to do is

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

$$a^{[2]} = \sigma(w^{[2]} a^{[1]} + b^{[2]})$$

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

for $i = 1$ to m :

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

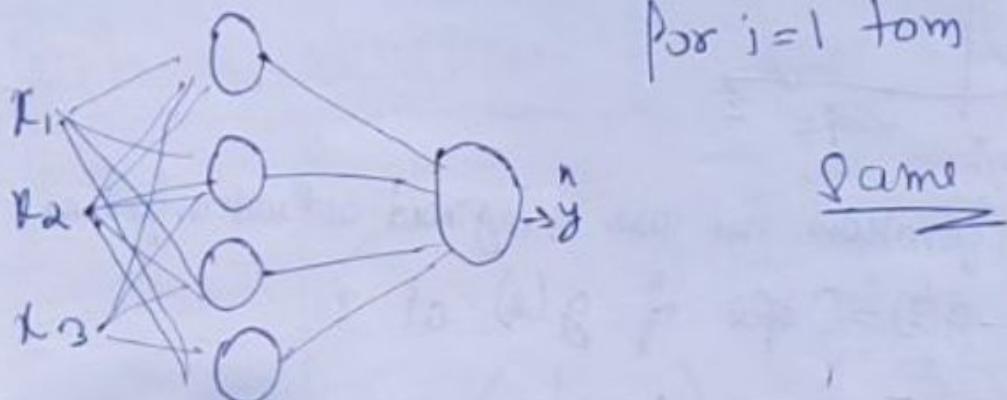
$$a^{[1]}(i) = \sigma(z^{[1]}(i))$$

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

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

m = drawing examples

④ EXPLANATION FOR VECTORIZED IMPLEMENTATION



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

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

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

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

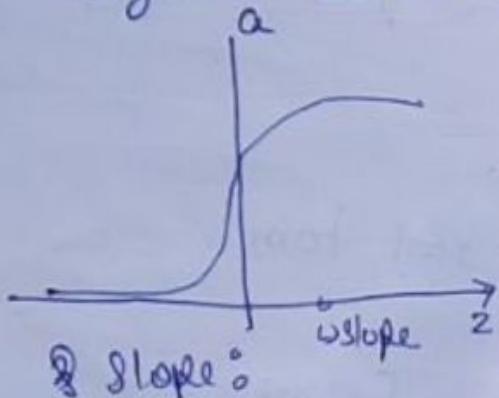
$$A^{[2]} = \sigma(z^{[2]})$$

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

⑤ ACTIVATION FUNCTIONS

In this we'll see the choices of activation function and how to find the slope of it.

① Sigmoid activation function



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

Slope of function in case of Sigmoid activation function

$$\begin{aligned} \text{slope of } g(z) &= \text{slope of } g(u) \text{ at } z \\ &= \frac{1}{1 + e^z} \left(1 - \frac{1}{1 + e^z} \right) \\ &= g(z)(1 - g(z)) \end{aligned}$$

* If z is very large $z=10$

$$z=10 \quad \cancel{\text{if}} \quad g(z) \approx 1 \quad (\text{g}(z) \text{ close to 1})$$

$$\frac{d}{dz} g(z) \approx 1(1-1) \approx 0 \quad (\text{very close to 0})$$

* If $z=-10$

$$\text{When } g(z) \approx 0 \quad (\text{close to 0})$$

$$\frac{d}{dz} g(z) \approx 0$$

$$\approx 0, (1-0) \approx 0$$

* If $z=0$

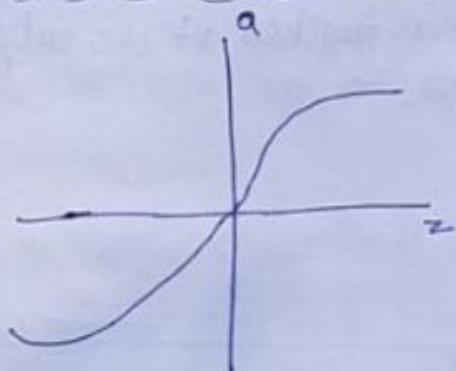
$$g(z) = \frac{1}{2}$$

$$\frac{d}{dt} g(w) = \frac{1}{2} (1 - \frac{1}{2}) = \frac{1}{4}$$

Note: A h calculus instead of writing $\frac{d}{dt} g(z)$ the shorthand for the derivative is $g'(z)$ which is written as $g'(t) = \frac{d}{dt} g(t)$

An neural network: $o = g(t) = \frac{1}{1+e^{-z}}$

Tanh activation function



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

$$\frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at particular point at } z \\ = 1 - (\tanh(z))^2$$

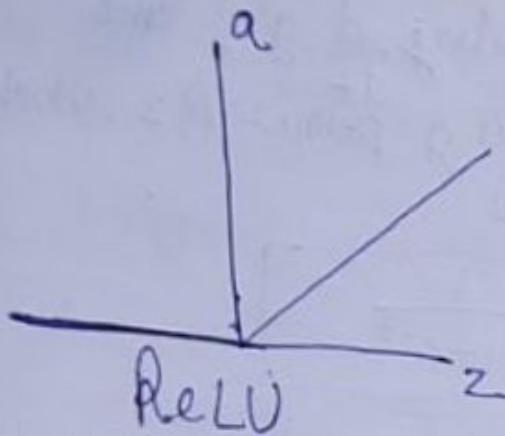
- * If $z=10$ $\tanh(z) \approx 1$
 $g'(z) \approx 0$

- * $z=-10$ $\tanh(z) \approx -1$
 $g'(z) \approx 0$

- * $z=0$ $\tanh(z)=0$
 $g'(z)=1$

Note: $a=g^{(1)}$, $g'(z) = 1-a^2$

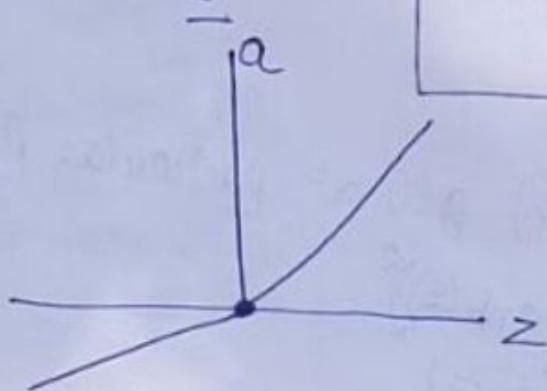
ReLU and Leaky ReLU



$$g(z) = \text{Max}(0, z)$$

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

Leaky ReLU



$$g(z) = \max(0.001z, z)$$

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

⑥ Why do we need non-linear functions?

If we use linear activation function on output layer, it will compute the output as a linear function of input features. It will be called as: $A = Wx + b$

also the output will be equal to z that is $A = z$

using linear activation is pointless. The composition of linear functions is still a linear function.

Remaining 6 is Here

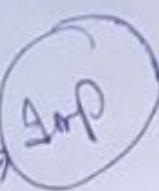
There is only one scenario where we tend to use a linear activation like price of house.

If we use a sigmoid function the output will range from (0) But the price will be more than 1\$. In this case we'll use linear function at the output layer.

GRADIENT DESCENT FOR NN

with one hidden layer

Note: ⑦ is on last page



NN. Parameters (with single hidden layer)

~~no. of input layers~~

$n^{[0]} = n_X$

~~no. of hidden layers~~

$n^{[1]} = \text{No. of hidden layer}$

~~no. of output layers~~

$n^{[2]} = \text{No. of Output Neurons} = 1$

$w_1 = \text{Shape is } (n^{[1]}, n^{[0]})$

$b_1 = \text{Shape as } (n^{[1]}, 1)$

$w_2 = (n^{[2]}, n^{[1]})$, $b_2 = \text{Shape } (n^{[2]}, 1)$

Cost function: $J(w^{[0]}, b^{[0]}, w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$

\hat{y}_i
predict

To draw the parameters & algorithm you need to perform gradient descent
gradient descent steps can be summarized as

- ① Compute predictions (\hat{y}_i , $i=1..m$)
 - ② Get derivatives: $dw^{[1]}, db^{[1]}, dw^{[2]}, db^{[2]}$
 - ③ Update $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$
- $w^{[1]} = w^{[1]} - \alpha * dw^{[1]}$, $b^{[1]} = b^{[1]} - \alpha * db^{[1]}$
- $w^{[2]} = w^{[2]} - \alpha * dw^{[2]}$, $b^{[2]} = b^{[2]} - \alpha * db^{[2]}$
- use partial derivatives here

lets look at the forward & backward propagation

Forward Propagation

$$z^{[0]} = w^{[0]} * A^{[0]} + b^{[0]} \quad \# A^{[0]} \text{ is } X$$

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

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

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

Backward Propagation

$$dz^{[2]} = A^{[2]} - y \quad \# y = [y^{(0)}, y^{(1)}, \dots, y^{(m)}]$$

$$dw^{[2]} = (dz^{[2]} * A^{[1]})^T / m$$

$$db^{[2]} = \text{Sum}(dz^{[2]}) / m$$

$$dz^{[1]} = (w^{[2]}.T + dz^{[2]} * g^{[1]})$$

element wise product

$$dw^{[1]} = (dz^{[1]} * A^{[0]})^T / m \quad \# A^{[0]} = X$$

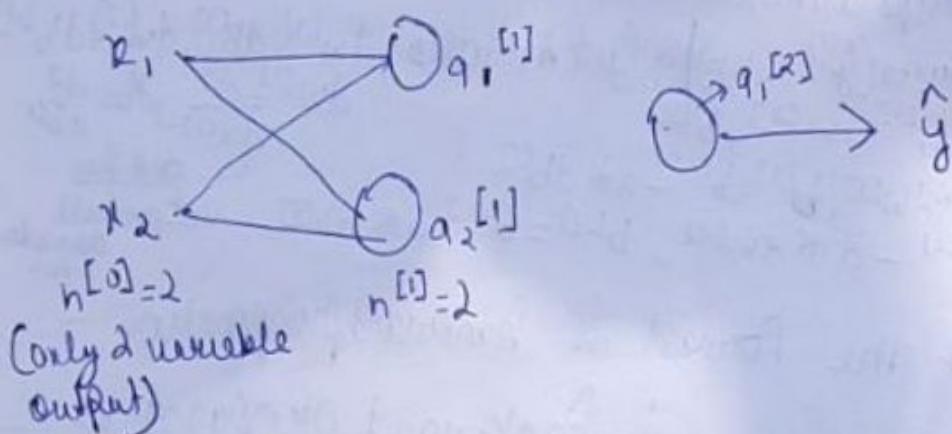
$$db^{[1]} = \text{Sum}(dz^{[1]}) / m$$

Watch video for proper explanation

⑨ RANDOM INITIALIZATION

When you change your Neural Network, it's important to initialize the weight randomly. For logistic regression, it was okay to initialize the weight to zero.

What happens if you initialize weights to zero



Suppose we initialise w and b to 0.

$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]}$$

If we initialise w to 0 then $a_1^{[1]} = a_2^{[1]}$

$$d_{z_1^{[1]}} = d_{z_2^{[1]}}$$

If we initialise w at 0 then $d_{z_1^{[1]}} = d_{z_2^{[1]}}$

The solution for to this is to initialize your parameter randomly.

If we initialize w and b at 0 then both nodes will give you same output, and all of the nodes will be symmetric.

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

↓
Multiplication by small no.

$$b^{[1]} = \text{np.zeros}(2, 1)$$

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

$$b^{[1]} = \dots$$

But why don't we multiply this with 100 or 200?
Because this will make activation functions more complex.

WEEK 4

~~What are deep NN?~~

① Deep L-layer NN.

* we can have any number of hidden layer in our model

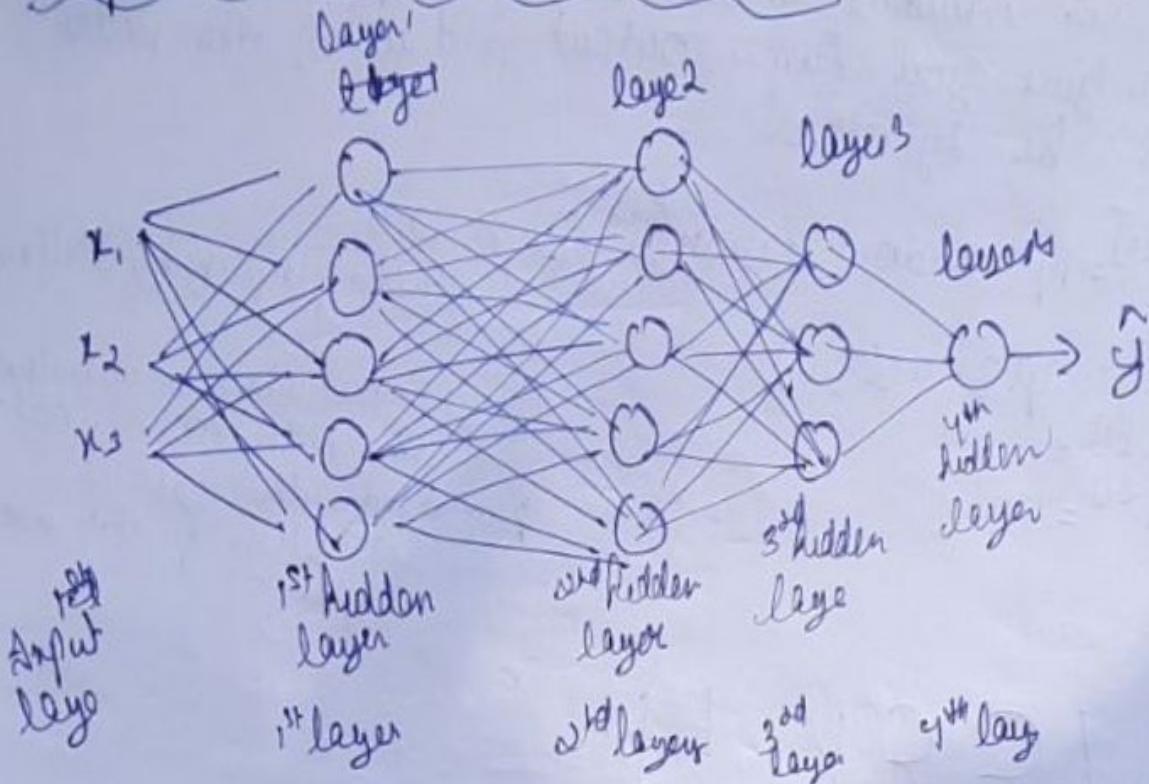
~~when we have 2 hidden layers we call it~~

* when we don't have any hidden layer in our model it is known as shallow NN.

* when we have 1 hidden layer it is called 2 layer NN (hidden layer + output layer)

* when we have many hidden layers it is called deep NN

Deep Neural Network notation

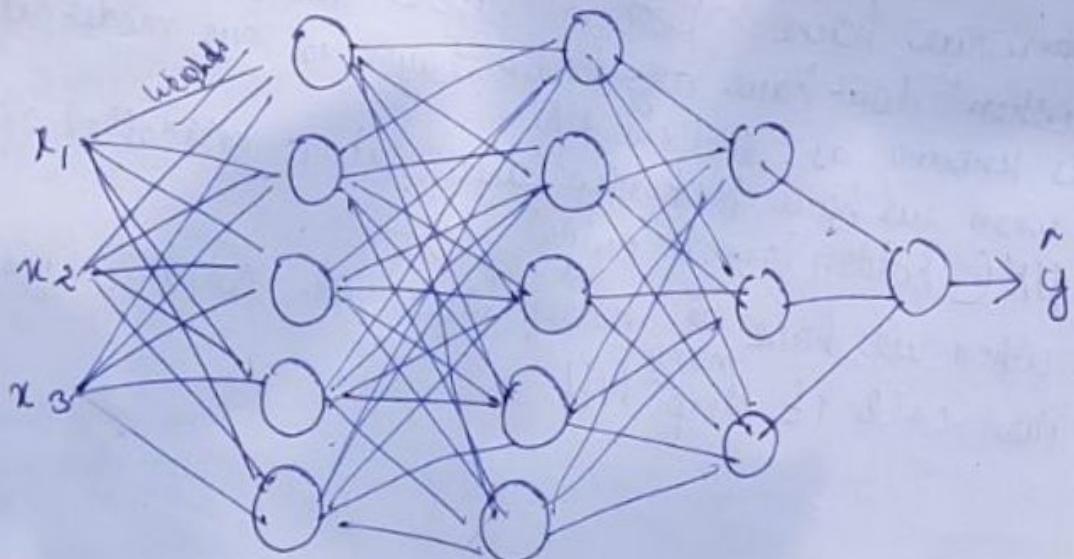


$$\text{No. of layers} = 4 \ (\# \text{ layers})$$

$$n^{[l]} = \# \text{ number of units in layer } l \quad n^{[1]} = 3, n^{[2]} = 5, n^{[3]} = 3, n^{[4]} = 1, 0$$

$a^{[l]}$: activation in layer

② FORWARD PROPAGATION IN A DEEP NETWORK



$$x : z^{[1]} = w^{[1]}x + b^{[1]} \quad \boxed{\text{for first layer}}$$

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

$$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]} \quad \boxed{\text{for layer 2}}$$

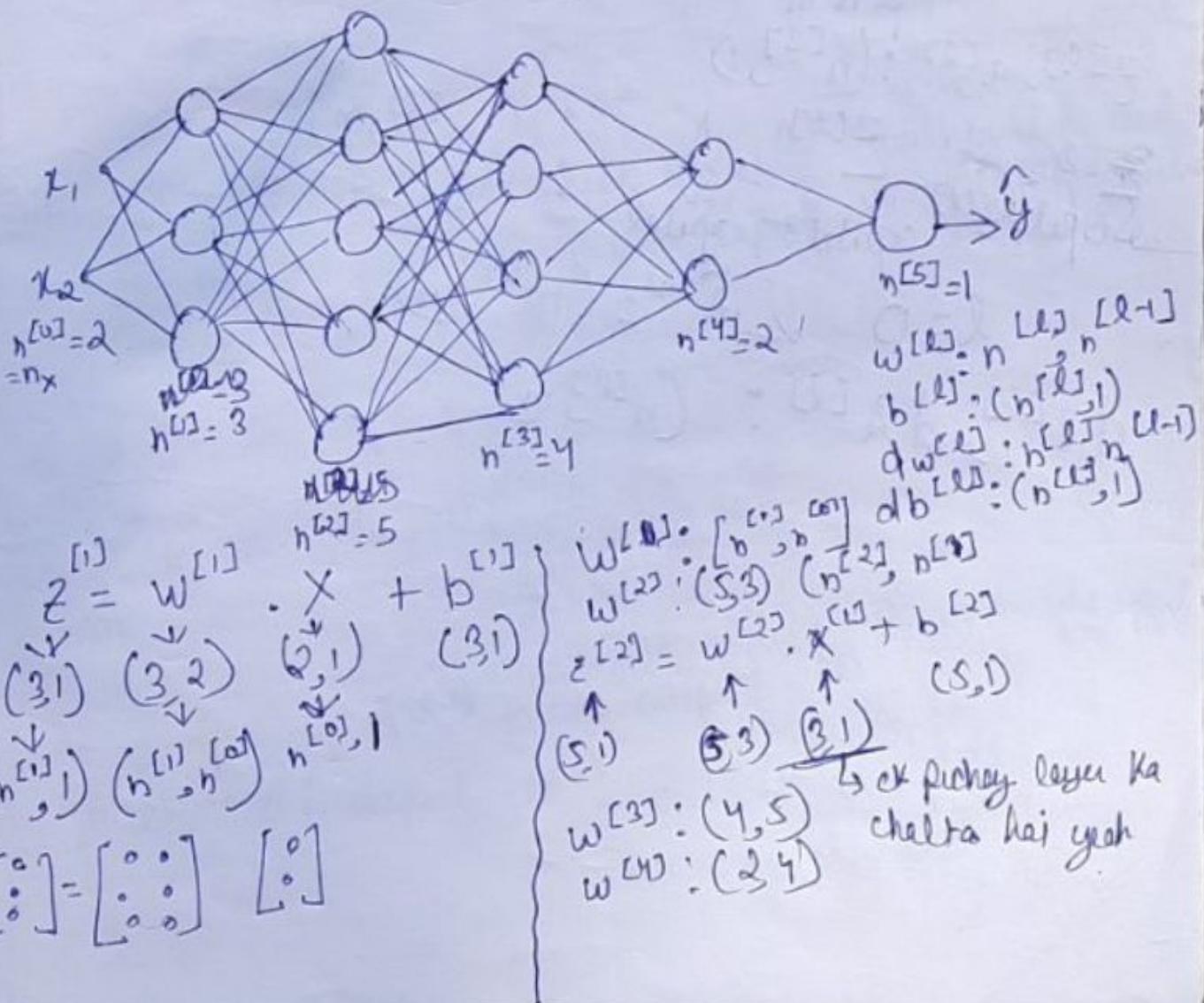
$$a^{[2]} = g^{[2]}(z^{[2]})$$

\vdots This we have to do for every layer

③ GETTING YOUR MATRIX DIMENSIONS

RIGHT

$l = \text{no. of layers} = 5$



Vectorized implementation

Name diagram

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

$$(n^{[1]}, 1) \xrightarrow{w^{[1]}, n^{[0]}} n^{[0], 1}) \rightarrow [n^{[0]}, 1]$$

$$\begin{cases} z^{[0]} = w^{[1]} \cdot x + b^{[1]} \\ (n^{[1]}, m) \xrightarrow{w^{[1]}, n^{[0]}} (n^{[0]}, m) \end{cases} \rightarrow (n^{[1]}, m)$$

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

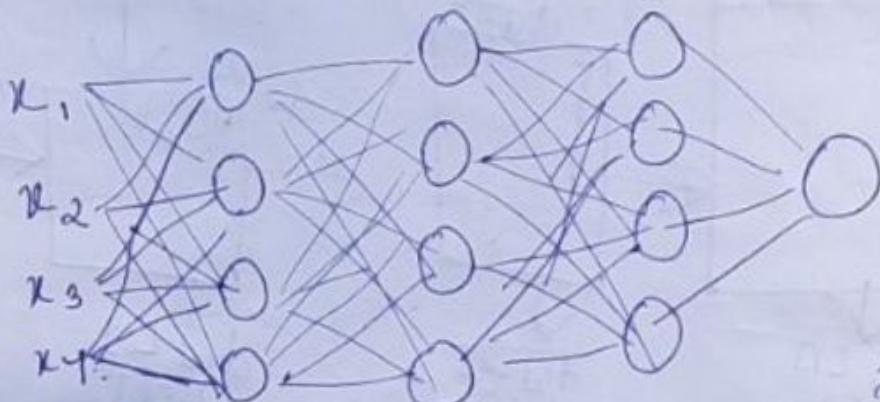
$$\begin{array}{ccc} z^{[l]} & \xrightarrow{\text{compute } z^{[l]}} & A^{[l]} : (n^{[l]}, m) \\ \downarrow & \downarrow & \downarrow \\ \text{compute } z^{[l]} & \text{compute } A^{[l]} & A^{[l]} : (n^{[0]}, m). \end{array}$$

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

④ WHY DEEP REPRESENTATION

- ↳ used in speech recognition
- ↳ used in photo object detection
- ↳ used in audio files

⑤ BUILDING BLOCKS OF DEEP N.N.



In this we cache
the value of $z^{(l)}$
which will be useful
in back-propagation

for layer l (any layer)
you have some parameters $w^{[l]}, b^{[l]}$
For forward prop you will will input
Forward: $a^{[l-1]}$ output $a^{[l]}$

↳ activation
from previous
layer

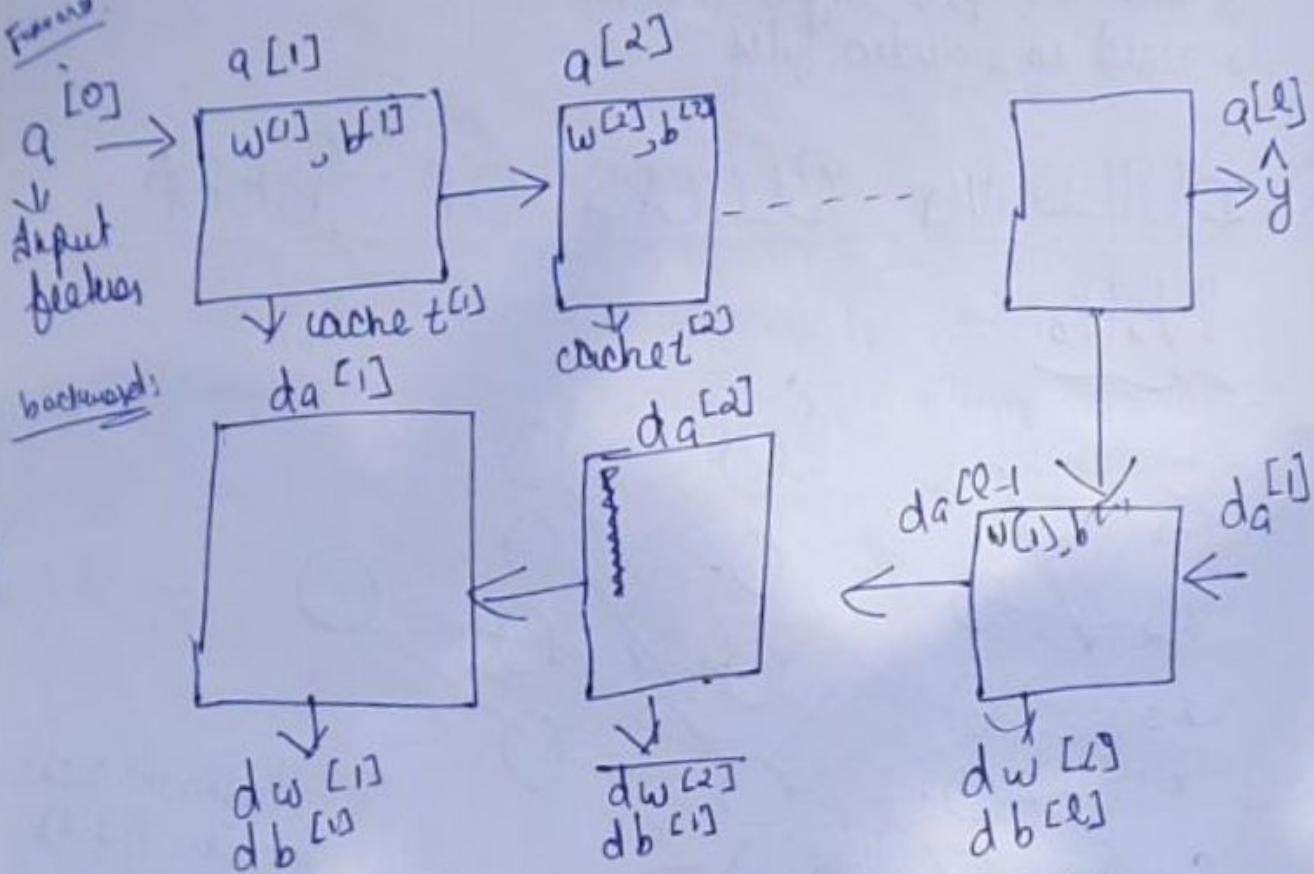
$$z^{(l)} = \underbrace{w^{(l)} a^{(l-1)}}_{\text{weights}} + b^{(l)}$$

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

that's how you go from the input
 $a^{(l-1)}$ to the output $a^{(l)}$

Backward: Input $d^{(l)}$, output $\frac{da^{[l]}}{dw^{[l]}}$

Forward and backward functions



⑥ FORWARD AND BACKWARD FUNCTIONS

Caches = values are saved of w and b in this for back propagation

it just save value we can use to save any value like variable

forward propagation for layer l

- input $a^{[l-1]}$
- output of $z^{[l]}$, as cache ($z^{[l]}$) → it include $w^{[l]}, b^{[l]}$ as well

$$g^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

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

and if you want vectorize implementation of this.

$$z^{[l]} = w^{[l]} \cdot A^{[l-1]} + b^{[l]}$$

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

Backward Propagation for layer 2

→ Input $d_a^{[l]}$

→ Output $d_a^{[l-1]}, d_w^{[l]}, d_b^{[l]}$

$$dz^{[l]} = d_a^{[l]} * g'(z^{[l]})$$

$$d_w^{[l]} = dz^{[l]} \cdot (A^{[l-1]})^T$$

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

$$d_a^{[l-1]} = w^{[l]T} * dz^{[l]}$$

Vectorize version

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

$$d_w^{[l]} = \frac{1}{m} \sum_{i=1}^m dz^{[l]}_i \cdot A^{[l-1]T}$$

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

$$dA^{[l-1]} = w^{[l]T} \cdot dz^{[l]}$$

⑦ PARAMETERS VS HYPERPARAMETERS

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, w^{[3]}, b^{[3]} \dots$

Hyperparameters: learning rate (α)
Iterations

hidden layer L

hidden units $n^{[1]}, n^{[2]}$

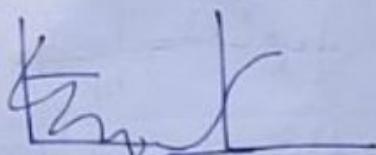
choice of activation function



#WEEK 3

⑦ DERIVATIVES OF ACTIVATION FUNCTION

For Sigmoid activation function.



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

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

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

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

This can also, a Hint! can also written as a $(1-a)$
written as

$$g'(z)$$

Same best for Sigmoid about this expression makes
sense?

If $z=10$ $g(z) \approx 1$ (\approx close to)

$$\frac{d}{dz} g(z) \Big|_{z=10} \approx 0$$

If $z=-10$ $g(z) \approx 0$

$$\frac{d}{dz} g(z) \Big|_{z=-10} \approx 0. (1-0) \approx 0$$

If $z=0$ $\nabla g(z) = \frac{1}{2}$

$$\frac{d}{dz} g(z) = \frac{1}{2}(1-\frac{1}{2}) = \frac{1}{4}$$

For Tanh function

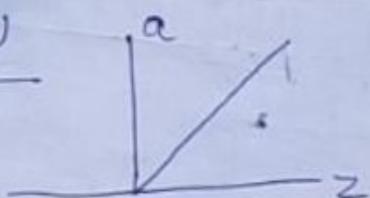
$$\cancel{\text{for}} \quad g(z) = \tanh(z)$$

$$g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z \\ = 1 - (\tanh(z))^2 \rightarrow \text{can also be written as} \\ 1 - \alpha^2$$

$$\begin{array}{c|c|c} \text{If } z=10 \quad \tanh \approx 1 & \text{If } z=-10 \quad \tanh \approx -1 & \text{If } z=0 \quad \tanh = 0 \\ g'(z) \approx 0 & g'(z) \approx 0 & g'(z) = 1 \end{array}$$

ReLU and Leaky ReLU

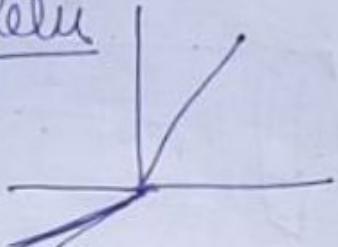
ReLU



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

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

Leaky ReLU



$$g(z) = \max(0.001z, z)$$

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

(3)

#2 IMPROVING DEEP N.N.

HYPERPARAMETER TUNING, REGULARIZATION AND OPTIMIZATION

WEEK 1

① Train / Dev / Test sets

↓ ↓ ↓
Training set development set (for validation) test dataset

Applied ML is a highly interactive process

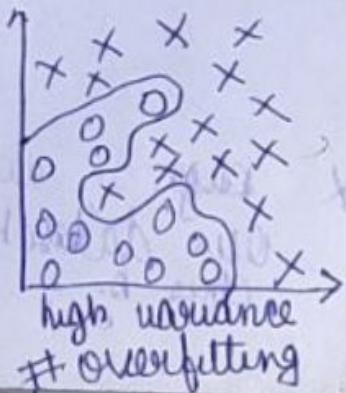
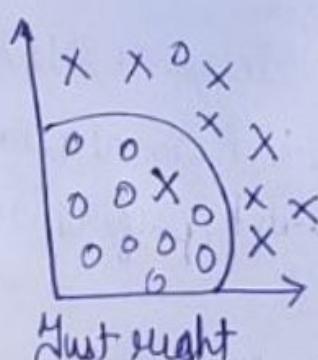
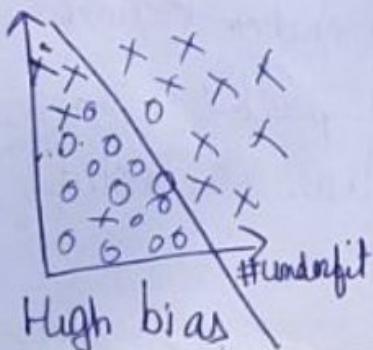
we can change:

- ① layers
- ② learning rate
- ③ activation function
and many more

Note: the testing data must be from the same dataset

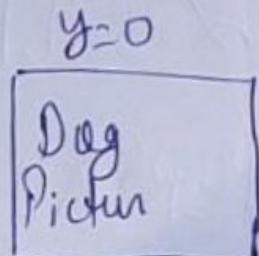
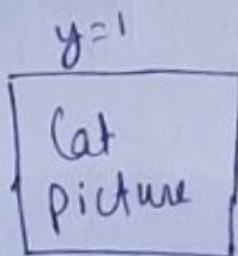
② Bias / Variance

In this we'll learn Bias and Variance tradeoff



Bias and Variance

Cat classification



Train set error	1%	15%	15%	0.5%
Dev set error	11%	16%	30%	1%
	high variance	high bias	high bias and high variance	low bias and low variance

- * When our model overfits some parts of the data like this



When it has Both Bias & Variance

③ BASIC RECIPE FOR M.L.

- * If we have high bias in our data (training data performance) we can solve this by-
 - ① - Bigger Network
 - ② Train longer
 - ③ Or trying different NN. architecture
- * ~~If~~ If we have high variance problem in our dataset (dev-set performance) we can solve this by -
 - with regularization
 - with dropout

- ① More data
- ② Regularization
- ③ NN. architecture structure.

Q If we want to decrease Bias in our data this will leads to little increase in Variance if we want to decrease Variance in our data this will leads to little increase in Bias this is called tradeoff b/w bias & variance. We have to create a proper balance b/w bias & variance.

④ REGULARIZATION → If our NN has overfitting problem (high variance) then we should try regularization

Let's take the example of logistic regression. We try to minimize the loss function

$$\min_{w,b} J(w, b)$$

This is how our cost function will look like if we add regularization.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m l d(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

↳ norm of w squared

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

w norm is also known as Frobenius norm

~~If we add $\|w\|_1$ in our cost function it is called L₁ regularization and if we add $\|b\|_1$ in our cost function it is called L₁ regularization~~

~~We generally use L₂ regularization because it is related to weights and weights has higher dimension than bias. So generally we use L₂ regularization~~

~~Note? We can also use L₁ & L₂ regularization together~~

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m l d(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2 + \frac{\lambda}{2m} \|b\|^2$$

Note: we call this $\|w\|^2$ as norm of w squared

So, why don't you regularize just the parameter w?
Why don't we add b as well like this.

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

because w is high dimension parameter and b is not
that will give us more difference

* We can also use L1 regularization so what is L1 regularization?

L1 regularization = when we add $\frac{\lambda}{2m} \sum_{j=1}^m |w_j| =$

$\frac{\lambda}{2m} \|w\|_1$ in our cost function. when we use L1 regularization this will has lots of zeros which will helps in compressing our model

So, L2 is used generally

& λ (lambda) in regularization is called regularization parameter. and lambda is another hyper parameter that you might have to tune

lets add a regularization term to our cost function
This is how it will look like now

$$J(w^{(0)}, b^{(0)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|^2$$

$$\|w^{(l)}\|^2 = \sum_{i=1}^{n^{(l)}} \sum_{j=1}^{n^{(l+1)}} (w^{(l)}_{ij})^2 \quad w \in (n^{(l)}, n^{(l+1)})$$

let's see how regularization works for a gradient descent algorithm:

previously we complete $d w$ using backprop

$$d w = (\text{from backprop}) \quad [\text{It gives the value of } \frac{d J}{d w^{[l]}}]$$

$$w^{[l]} = w^{[l]} - \alpha d w^{[l]}$$

But this is before we added the extra regularization

$$d w^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]}$$

$$w^{[l]} = w^{[l]} - \alpha d w^{[l]} + \frac{\lambda}{m} w^{[l]}$$

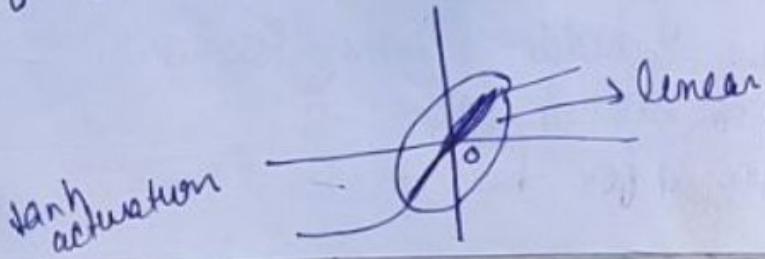
L2 regularization is also called weight decay.

$$w^{[l]} = w^{[l]} - \alpha [(\text{from backprop}) + \frac{\lambda}{m} w^{[l]}]$$

$$= w^{[l]} - \frac{\lambda}{m} w^{[l]} - \alpha (\text{from backprop})$$

⑤ WHY REGULARIZATION REDUCE OVERFITTING

If the regularization becomes very large, the parameter w very small, so z will be relatively small (ignore b for now). So z will take small range of value and so the activation function if is is tanh will be relatively linear and your whole N.N. will be computing something not too far from a big linear function, which are much less useful.

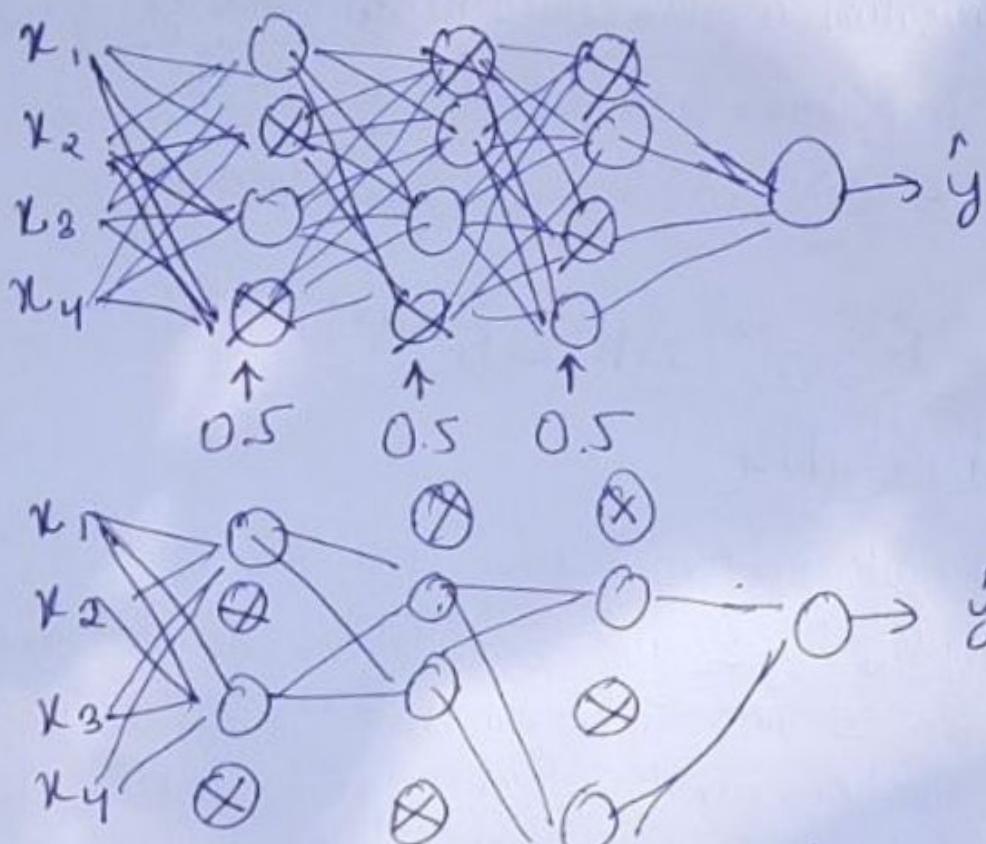


because these values of $z^{[l]}$ will be near the origin

⑥ DROPOUT REGULARIZATION

This is a very powerful regularization technique called dropout.

Suppose this neural network is overfitting on the training data. Suppose we add a dropout of 0.5 to all these units. The model will randomly remove 50% of the units from each layer and we finally end up with a much simpler network.



This has power to be a very effective regularization technique. How can we implement lets check

- * Suppose we have 3 hidden layers. For now, we will consider the third layer, $l=3$.
The dropout vector d for third hidden layer

can be written as:

d^3 = dropout vector for the layer 3

$d^3 = \text{np.random.rand}(a^3.\text{shape}[0], a^3.\text{shape}[1]) < \text{keep_prob}$

Here, keep_prob is the probability of keeping a unit.

So, if keep_prob=0.8 then this means that there is a 0.2 chance of elimination an hidden unit.

$a_3 = a_3$ has the activation that you computed

$a_3 = \text{np.multiply}(a^3, d_3)$

This a^3 value will be reduced by a factor of $(1 - \text{keep_prob})$.
So to get the expected value of a^3 , we divide the value

$a^3 / = \text{Keep_dum}$

example:-

No. of units in the layer = 50

Keep - prob = 0.8

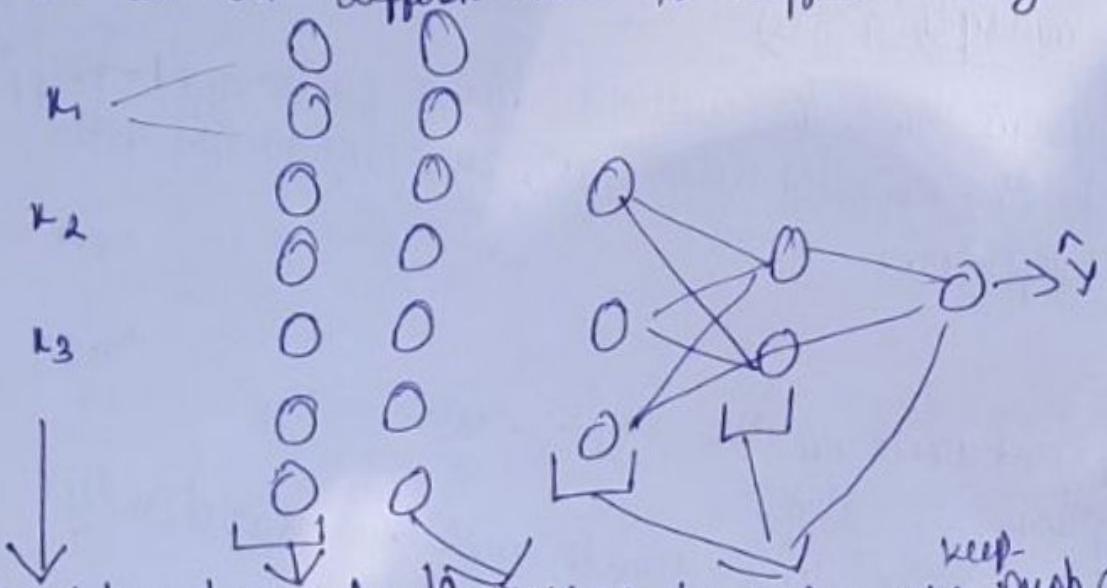
So, 20% of the total units will be randomly shut off.

Note that dropout is only done at the time of training the model (not during the test phase). The reason for doing this is because:

- We don't want our output to be random.
- Dropout adds noise to the prediction.

(7) UNDERSTANDING DROPOUT

- * When we use dropout it just knows knock neurons in a hidden layer randomly.
- * If we have more nodes in hidden layers, we can know more neuron using Dropout.
- * Following is the example of different knockoff value for different layer of N.N.
- * we can set different value for different layers.



for input layer, keep_prob should be higher because generally we don't want our data to drop out early.	we can have keep_prob higher or less because we have less nodes that means here suppose our keep_prob = 0.3 then drop out value = 0.3 30% of the nodes will drop out	Suppose we have less nodes because we keep_prob is 0.5 here that means that 50% of the nodes will drop out	we can use prob at also here no neuron will be dropped.
--	---	---	---

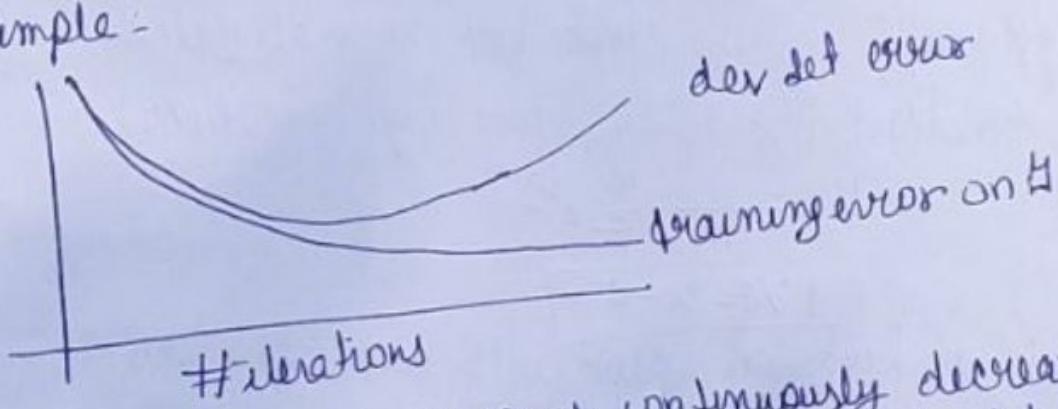
- * We don't use ~~Dropout~~ when our model is not overfitted (high bias or high variance).

⑧ OTHER REGULARIZATION METHOD

① Data Augmentation: Suppose we are building an image classifier model and are lacking data with the training data. In such cases, we can use data augmentation i.e. applying some changes in the image by flipping, random crops, rotating etc or adding distortions in the image etc. These can help in getting more training data and helps in reducing overfitting.

② Early Stopping:

Example -



Here, the training error is continuously decreasing with respect to time. On the other hand, the dev set error is decreasing initially before increasing after a few iterations. We can stop the learning at the point where the dev set error starts to increase. This is called early stopping.

~~VIDEO 2~~ Setting up your optimization function

~~① MINI-BATCH GRADIENT DESCENT~~

~~② NORMALIZING INPUTS~~ → It helps in speeding up our NN

Here we have 2 input features $x_1, 2x_2$



$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

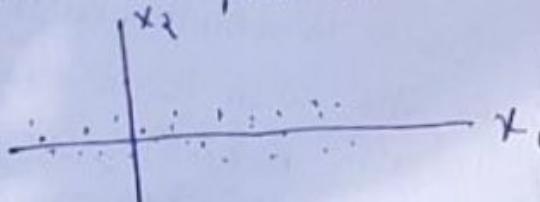
following are the steps for normalization.

1. Subtract the mean from input features:

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

$$x := x - \bar{x}$$

This is how our plot will look after this



→ Here x_1 has the large variance than x_2

2. Next we normalize the variance

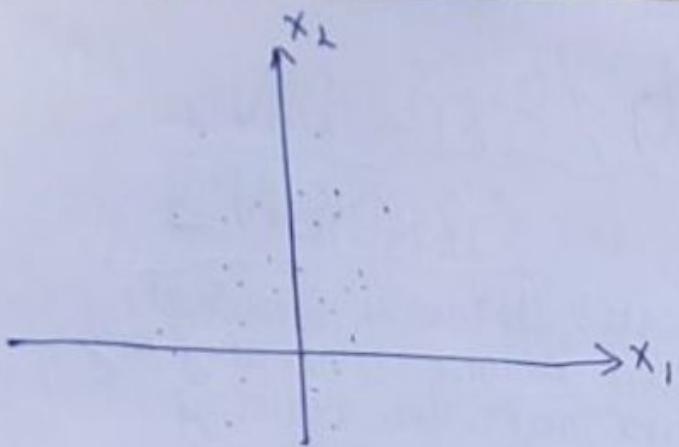
Here, feature x_1 has larger variance than x_2

$$\sigma^2 = \frac{1}{m} \sum_i (x_i - \bar{x})^2 \rightarrow \text{this is elementwise squaring}$$

\rightarrow square (thus multiplication)

we divide the feature with the variance

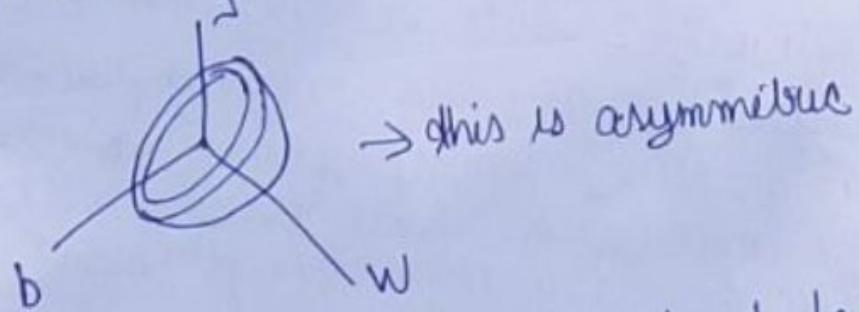
This is how our input looks like this after that



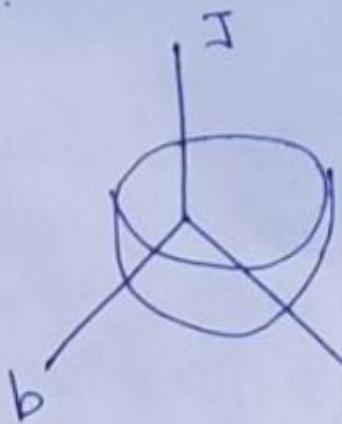
The same mean and variance to normalize the test set as well. ~~unrelated~~ because we want the same transformation to happen on both the train & test data.

How does normalizing the data make the algorithm faster?

An unnormalized data, the scale of features will vary and hence there will be variation in the parameter learnt from each feature. This will make the cost function asymmetric.



whereas, in the case of normalized data, the scale will be the same & cost function will also be symmetric.



Normalizing the inputs makes the cost function symmetric. This makes it easier for the gradient descent algorithm to find the global minima more quickly. This makes algorithm to run much faster.

⑩ VANISHING / EXPLODING GRADIENTS

while training deep Neural Networks, sometimes the derivatives (slopes) can become either very big or very small. It can make the learning phase quite difficult. This is the problem of vanishing/exploding gradients. Suppose a N.N. with 1-layer with 2 input and we initialize the large weights:

$$w^{[1]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

The final output of 1st layer will be

$$\hat{y} = w^{[1]} \begin{bmatrix} 1.5 \\ 0 \end{bmatrix}^{-1}$$

for deeper network, L will be large, making the gradient very large and the learning process slow. Similarly, using small weights make the gradients very small as a result learning will be slow. So how to deal with this problem?
How to initialize weights?

⑪ WEIGHT INITIALIZATION FOR DEEP NETWORK

It turns out that the partial solution to this doesn't solve it entirely but helps a lot, better or more useful choice of the random initialization for your Neural Network.

example:-

for sigmoid

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

larger \rightarrow smaller w_i

let $b=0$

$$\text{Var}(w_i) = \frac{1}{n}$$

Note: we always want our weight to be smaller because ~~of~~ some weight do lots of multiplication

$$w^{[1]} = \text{np.random.randn}(\text{shape}) * \text{np.sqrt}(\frac{1}{n^{(d-1)}})$$

In case of ReLU Activation function-

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

$$w^{[1]} = \text{np.random.randn}(\text{shape}) * \text{np.sqrt}(\frac{2}{n^{(d-1)}})$$

The primary reason behind initializing the weights randomly is to break symmetry. We want to make sure that different hidden units learns different patterns. There is one more technique that can help ensure that our implementation are current and will run quickly.

⑫ NUMERICAL APPROXIMATION OF GRADIENTS

When you implement back propagation you'll find that theres a test called gradient checking that can really help you make sure that your implementation of back prop is correct.







WEEK 2

① MINI-BATCH GRADIENT DESCENT

what happens if we have a large training set say $m = 5,000,000$? if we train our model with all training examples it will take lot of time. So Mini-Batch Gradient descent is the solution for this.

So, we split of data sets in mini-batches up to 1,000 examples each. It means for 1,000 examples we have to create 5,000 batches.

This is how our training data set will look:

$$X = [x^{(1)} \ x^{(2)} \ x^{(3)} \dots x^{(1,000)} | x^{(1,001)} \dots x^{(2,000)} | \dots x^{(m)}]$$

\downarrow

$x^{(1,000)}$
↳ 1st Batch

$$Y = [y^{(1)} \ y^{(2)} \ y^{(3)} \dots y^{(1,000)} | y^{(1,001)} \dots y^{(2,000)} | \dots y^{(m)}]$$

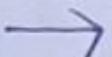
\downarrow

$y^{(1,000)}$
↳ 1st Batch
↳ input

Here $x_i^{(k)}, y_i^{(k)}$ represents the i th mini-batch input & output.

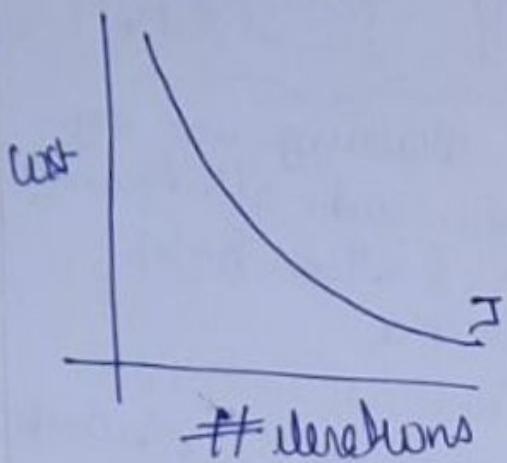
② UNDERSTANDING MINI-BATCH GRADIENT DESCENT

Process



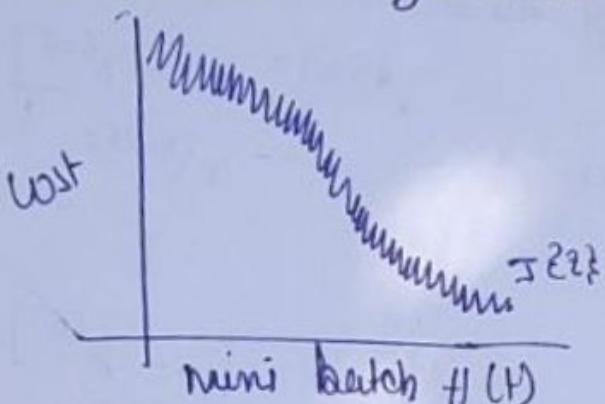
Training with mini batch gradient descent

Batch gradient descent



With batch gradient descent
on every iteration you
go through the entire
training set and you'd expect
the cost to go down on
every single iteration.

Mini-batch gradient descent



It needs to go downwards but
it is going to be noisy

On mini-batch gradient descent
if you plot progress on your
cost function, then it will not
decrease on every iteration.

If we plot cost function $J^{(t)}$
using $x^{(t)}, y^{(t)}$ (This is how our
plot will look like)

Note: in mini-batch gradient descent is noisier because maybe
other others are mislabeled in 1st batch and early fit on second

Choosing your mini-batch size

① If your minibatch size = m : Then it is called batch
minibatch = size of training rows Gradient descent
 $(x^{(1)}, y^{(1)}) = (x, y)$

② If your mini-batch size = 1 : This is called
each training example is its own
mini-batch Stochastic gradient
descent

It can be extremely noisy and takes more time to reach global minima.

3. If the mini-batch size is between 1 to m :

It is mini-batch gradient descent. The size of the mini-batch should not be too large or too small.

Below are a few general guidelines to keep in mind while deciding the mini-batch size:

1. If the training set is small, we can choose a mini-batch size size $\delta m \leq 1000$.

2. For a larger training set, typical mini-batch size are 64, 128, 256, 512
↳ These are defined by power of 2 for example $2^6, 2^7, 2^8, 2^9$ and it is to make sure that the mini-batch size fits your GPU / CPU memory not fit

③ EXPONENTIALITY WEIGHTED AVERAGE

Following are the daily temperature of London 2011

$\theta_1 = 40^\circ\text{F}$ (1st day of year)

$\theta_2 = 49^\circ\text{F}$

$\theta_3 = 45^\circ\text{F}$

=

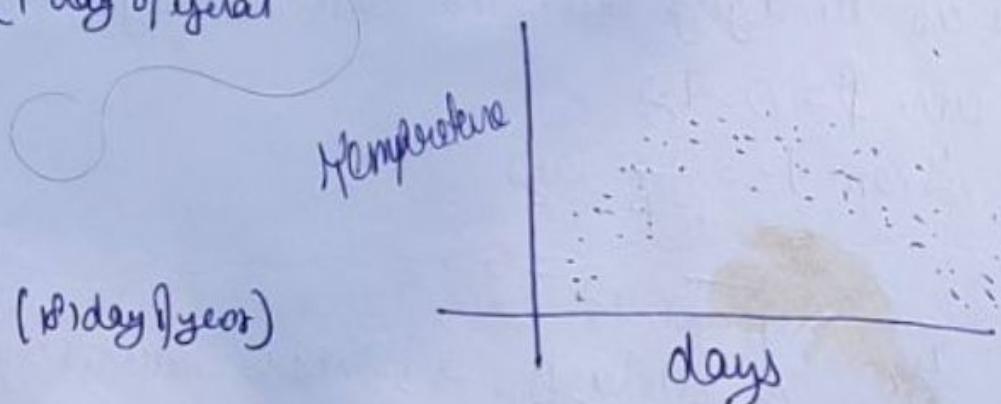
$\theta_{180} = 60^\circ\text{F}$

$\theta_{181} = 56^\circ\text{F}$ (1st day of year)

=

$\theta_{365} = ?$

Exponentially weighted averages are used by many optimisation algorithms



when we plot this we'll get this graph,

If you want to compute the moving average or a moving average of the temperature

lets initialize $V_0 = 0$

Then on every day $V_t = 0.9V_{t-1} + 0.1D_t$ [$D_t = \text{the day temperature}$]

$$V_1 = 0.9V_0 + 0.1D_1 \quad (\text{1st day temperature})$$

$$V_2 = 0.9V_1 + 0.1D_2 \quad (\text{2nd day temperature})$$

If we plot this we'll get



You'll get a moving average of what's called an exponentially weighted average of the daily temperature

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

$\beta = 0.9$

temperature of the day
from the previous day.

V_t = approximately average over $\frac{1}{1-\beta}$ days temperature
So, for example when Beta goes 0.9 you could think of this as averaging over the last 10 days temperature.

* If our $\beta = 0.98$

then $V_t = \frac{1}{1-\beta}$ days

$$= \frac{1}{1-0.98} = 50$$

or $\times 50$ days This is averaging over roughly the last 50 days temperature

a high value of Beta the plot you'll get is much smoother because you're now averaging over more days of temperature

④ Day UNDERSTANDING WEIGHTED AVERAGE

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

$$V_{1,00} = 0.9 V_{99} + 0.1 Q_{99}$$

$$V_{98} = 0.9 V_{97} + 0.1 Q_{98}$$

$$V_{1,00} = 0.1 Q_{1,00} + 0.9 \xrightarrow{\text{again we form}} (0.0 Q_{99} + 0.9 Q_{98})$$

This is how $V_{1,00}$ is calculated

⑤ BIA SCORRECTION IN EXPONENTIALLY WEIGHTED AVERAGE

There is one technical detail called bias correction that can make your computation of these averages more accurately.

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

When you're implementing a moving average you initialise it with $V_0 = 0$

$$V_0 = 0$$

$$V_1 = 0.98 V_0 + 0.02 Q_1 (\text{if } Q_1 = 40)$$

$$V_1 = \beta (V_0 = 0)$$

It's not a good estimate of the first day temperature

$$V_0 = 0.98 V_0 + 0.02 Q_0$$

= —

It is not a good estimate for the initial phase of your estimate.

So, here instead of using the previous equation i.e.,

$$V_t = \beta * V_{(t-1)} + (1 - \beta) * \theta_t$$

we include a bias correction term:

$$V_t = [\beta * V_{(t-1)} + (1 - \beta) * \theta_t] / (1 - \beta_t)$$

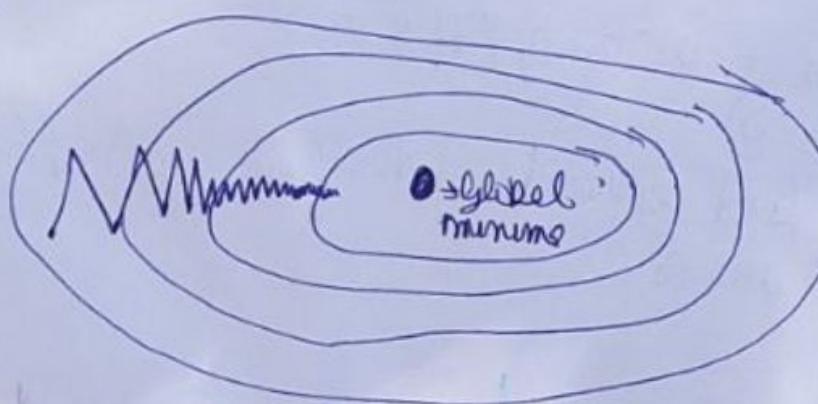
When t is small, β_t will be large, resulting in a smaller value of $(1 - \beta_t)$. This will make the V_t value large which ensures that our predictions are accurate.

⑥ GRADIENT DESCENT WITH MOMENTUM

There's an algorithm called momentum, or gradient descent with momentum that almost always works faster than the standard gradient descent algorithm.

In one sentence, the basic idea is to compute an exponentially weighted average of your gradients, and then use that gradient to update your weights instead.

The idea of gradient descent with momentum is to calculate the exponential weighted average of gradients and use them to do update weight. Suppose we have a cost function look like this:



Now we want to reach global minima and we want to reach that point.

One more way could be to use a larger learning rate. But that would result in large update steps, and we might not reach global minima. Also too small a learning rate makes the gradient descent slower.

* We want a slower learning in the vertical direction and a faster learning in the horizontal direction which will help us to reach the global minima much faster.

Compute Δw , Δb on current mini batch using momentum -

$$\nabla \Delta w = \beta * \nabla \Delta w_t + (1 - \beta) * \nabla \Delta w$$

$$\nabla \Delta b = \beta * \nabla \Delta b_t + (1 - \beta) * \nabla \Delta b$$

Update Weights

$$w = w - \alpha * \nabla \Delta w$$

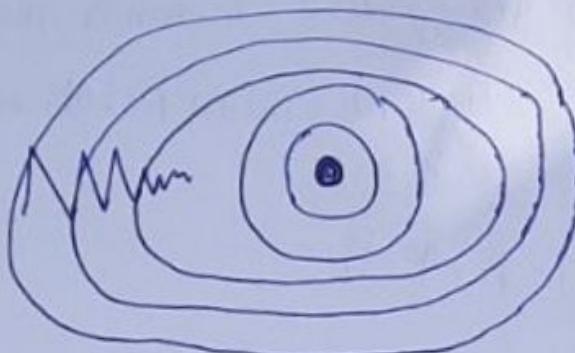
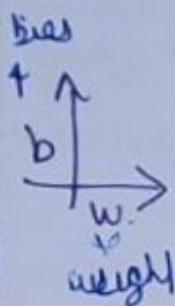
$$b = b - \alpha * \nabla \Delta b$$

(Add)

Here we have two hyperparameters i.e., α & β . The role of Δw and Δb in the above equation is to provide momentum. $\nabla \Delta w$ and $\nabla \Delta b$ provides velocity, and β acts as friction and prevents speeding over the limit. Consider a ball rolling down - $\nabla \Delta w$ and $\nabla \Delta b$ provide velocity to the ball and make it move faster. We do not want our ball to speed up so much that it misses the global minima and hence β acts as friction.

(7) RMSPROP

RMSPROP means root mean square prop that can also speed up gradient descent.



The vertical axis is the parameter b and horizontal axis is the parameter w .

- * and we want to slow down the learning in b direction or vertical direction and speed up the ~~is~~ horizontal direction.

On iteration t :

compute d_w, d_b on current mini-batch this is element wise squaring operation

$$S_{dw} = \beta_a * S_{dw} + (1 - \beta_a) * d_w^2 \leftarrow \text{small} \rightarrow \cancel{\text{large}}$$

$$V_{db} = \beta_a * V_{db} + (1 - \beta_a) * d_b^2 \leftarrow \text{large}$$

Update weights S_{dw} is relatively small than V_{db} so divided by small number

$$w = w - \alpha * (d_w / \sqrt{S_{dw}})$$

$$b = b - \alpha * (d_b / \sqrt{V_{db}})$$

when vertical movement is done by b (bias) & horizontal by w (bias)

Here we are dividing w with a small number and b with a large number to boost horizontally

Vertical movement is smaller than horizontal

The slope in the vertical direction (d_b) is steeper, resulting in a large value of S_{db} . As we want slow learning in the vertical direction, dividing d_b with S_{db} in update step will result in smaller change in b . Hence, learning in the vertical direction will be less. Similarly, a small value of S_{dw} will result in faster learning in the horizontal direction.

⑧ ADAM OPTIMIZATION ALGORITHM

Adam is combination of momentum and RMSprop

There are a range of hyperparameters used in adam and some of the common ones are:

- o Learning rate α : needs to be tuned
- o Momentum term β_1 : common choice is 0.9
- o RMS prop term β_2 : common choice is 0.999

Adam helps to train a NN model much more quickly than the techniques we have seen earlier

To implement Adam you would initialize $V_{dw} = 0$, $S_{dw} = 0$ and similarly $V_{db} = 0$ & $S_{db} = 0$. And then on iteration i and T you would compute the derivatives: ~~compute dw, db~~

$$\textcircled{1} \quad V_{dm} = \beta_1 V_{dw} + (1 - \beta_1) dw, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

$$\textcircled{2} \quad S_{dw} = \beta_2 S_{dm} + (1 - \beta_2) dw^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$$

So, $\textcircled{1}$ is momentum like update with V_{dm} and $\textcircled{2}$ is RMSprop update with S_{dm} .

In the typical implementation of Adam, you do implement bias correction.

corrected means after bias correction

$$V_{dm}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1^t)$$

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

$$\omega_t = \omega - \alpha * \frac{V_{dm}^{\text{corrected}}}{\sqrt{S_{dm}^{\text{corrected}}} + \epsilon} \quad b_t = \beta - \alpha * \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

ϵ is very small number suppose if denominator become 0 then we'll don't get value of w & b . & that's why we ~~don't~~ use ϵ (Epsilon). This algorithm has a numbers of hyper parameters.

↳ Given in Slides

beta1 :- 0.9 is recommended

beta2 :- 0.999 is recommended

Epsilon : 10^{-8} is n

AM = Adaptive Moment Estimation

— (dw) — known as first movement.

— (d^2w) — known as second movement.

⑨ LEARNING RATE DECAY

If we slowly reduce the learning rate over time, we might speed up the learning process. This process is called learning rate decay.

1 epoch = 1 pass through the data

2 epoch = 2^{nd} pass through the data

$$\alpha = \frac{1}{1 + \text{decay_rate} * \text{epoch_num}} \alpha_0$$

↓ some initial learning rate α_0

learning rate

→ decay_rate has become another hyperparameter which you might need to tune.

If $\alpha_0 = 0.02$

decay_rate = 1

Then

Epoch	α
1	$\frac{1}{1 + 1 * 1} * \alpha_0 = 0.02$
2	0.67
3	0.5
4	0.4

Here learning rate is decreasing

There are also many formulas to calculate learning rate decay

10. LOCAL OPTIMA IN N.N.

It refers to when we ^{say to} find the minimum value of our cost function we get stuck in a local optima means stuck in near minimum value of it. But this is not the minimum value for our cost function

WEEK 3

④ HYPERPARAMETERS

① Tuning Process

following are the few hyperparameters:

- ① Learning rate - α (most important)
- ② Momentum - β
- ③ Adam's hyperparameter - B_1, B_2
- ④ Number of hidden layer
- ⑤ Learning rate decay
- ⑥ Mini-Batch size

② PICKING HYPERPARAMETERS

AT RANDOM

① Hidden layer

You should have 2-4 hidden layer in your model N.N.

② learning rate

α can be low as 0.00001 - - - to 1
(or any small number)

BATCH NORMALIZATION

① NORMALIZING ACTIVATIONS IN A NETWORK

In deep N.N. we have lots of layers (hidden layers) which has nodes \rightarrow means lots of activation

In this we normalize the mean and variance ~~in order~~ so that the weights of the next layer can be update faster

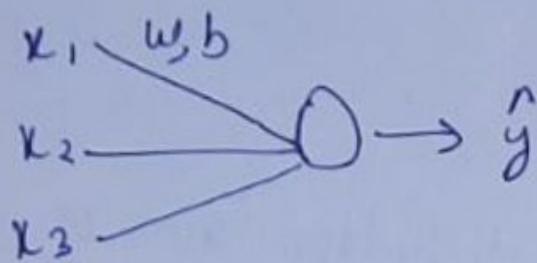
Quick Recap

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

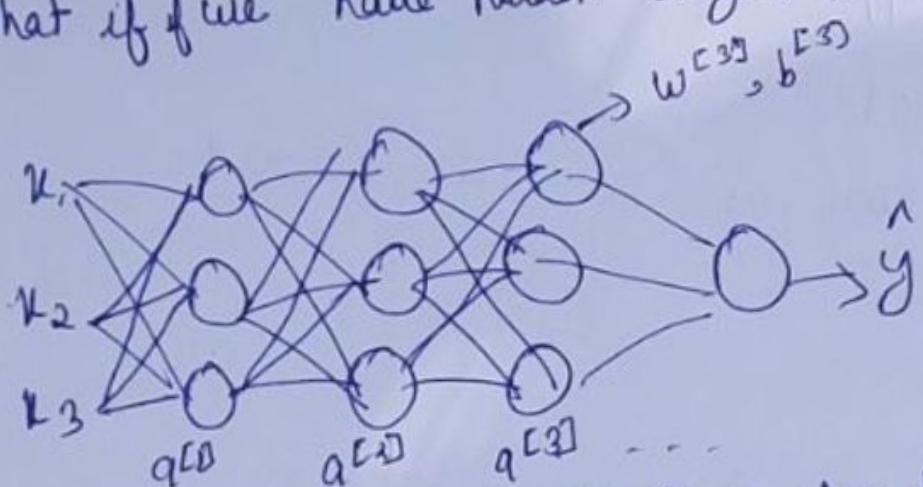
$$\sigma^2 = \frac{1}{m} \sum (x^{(i)} - \mu)^2 \rightarrow \text{elements by squaring}$$

$$X = X / \sigma^2$$

this work only with this type of model



what if we have hidden layers in our model.



So, we want to train the parameters w & b
and b_3 (3^{rd} layer w & b)

& we can normalize the mean and variance of $a^{(2)}$ to
make the training of $w^{(3)}$ & $b^{(3)}$ for
more efficient

By normalizing the $a^{(2)}$ layer w & b it will
make $a^{(3)}$, $b^{(3)}$ work faster

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

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

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

[Normalizing
value]

$$y^{(i)} = \gamma z_{\text{mean}}^{(i)} + \beta$$

Here γ & β are learnable parameters

② FIXING BATCH NORM INTO A N.N.

③ WHY DOES BATCH NORMALIZATION WORK?

- The first reason is the same reason why we normalize X .
- The second reason is that batch normalization reduces the problem of input values changing (shifting).
- Batch normalization does some regularization.
 - Each mini batch is scaled by mean/variance computed of that mini-batch.
 - It also adds some noise of the values to the values $Z[i]$ within that mini batch. So similar to dropout it adds some noise to each hidden layer activation.

FOOT 1000
FOOT 1000

- o This has a slight regularization effect. (4)
 - o Using bigger size of the mini-batch you are reducing noise and therefore regularization effect.
- ② Don't rely on batch normalization as a regularization. It's intended for normalization of hidden units, activation & therefore speeding up learning. For regularization use other regularization technique (L^2 or dropout)

④ ~~BATCH NORMALIZATION~~ TEST TIME

- o When we train a N.N. with batch normalization, we compute the mean & variance of the mini-batch.
- o In testing we might need to process examples at a time other mean & variance of one example won't make sense.

- We can use weighted average across the mini-batches
- We have to compute an estimated value of mean & variance to use it in testing time.
- ~~We can use the weighted average across the mini-batches~~
- We ~~can~~ can use the weighted the weighted average across the mini-batches
- We will use the estimated value of the mean & variance to test
- This method is called running average
- An ~~practice~~ ^{will use} most deep learning framework ~~and~~ and it will implement ~~etc.~~ these things by default

Ans:

7

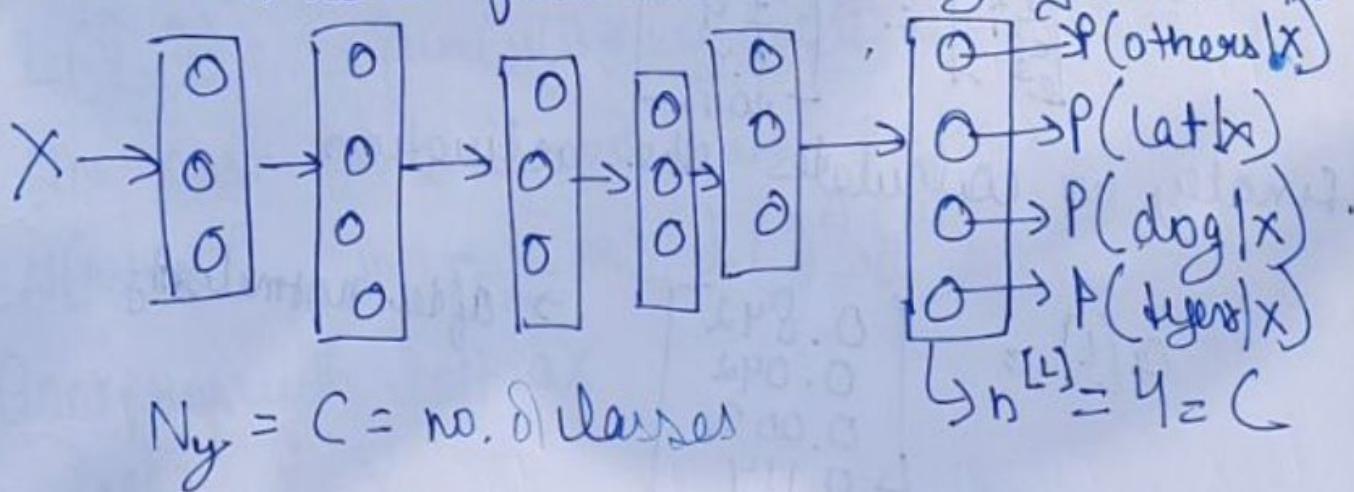
Ans :-

MULTI-CLASS CLASSIFICATION

① SOFTMAX REGRESSION

Binary classification means dealing with just 2 classes. But when we have more than one classification we use Softmax (multi class function). Suppose we have to recognise cat, dogs & tigers in set of images.

At the output layer, instead of having a single unit, we have units equal to the total number of class which is 4 in our example not 3 because one class will represent none of them option. Each unit tells us the probability of the image falling in different class. Since it tells the probability, the sum of values from unit is always equal to 1.



In the last layer we will have to activate the softmax activation function instead of the sigmoid activation.

Activation function $t^W = w^{[1]} a^{(l-1)} + b^{[1]}$

$t = e^{t^W} \rightarrow t$ will be $(4, 1)$ matrix
temporary variable

output $= a^{[L]} = \frac{e^{z^{[L]}}}{\sum_t}$ (it is also 4×1 vector)

let's understand this with an example.
Consider the output from the last hidden layer:

$$z^{[1]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

We then calculate t using the formula given above:

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

Finally we calculate the activation

$$a^{[L]} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \rightarrow \text{after normalizing}$$

② TRAINING A SOFTMAX CLASSIFIER

- There's an activation which is called hard max, which gets 1 for the maximum & zeros for others
- The softmax is a generalization of logistic activation function to C classes. If C=2 Softmax reduces to logistic regression
- The softmax is generalization of logistic activation function to C classes
- The Softmax name came from softening the values & not hardening them like Hard max
The loss function used with Softmax

$$L(\hat{y}, \hat{y}) = -\sum y(i) * \log(\hat{y}[i]) \quad \# i = 0 \text{ to } C-1$$

The cost function used with Softmax.

$$J(w^{(0)}, b^{(0)}, \dots) = -\frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, \hat{y}^{(i)})$$

Backprop with Softmax

$$d_2^{(cls)} = \hat{y} - y$$

derivative of Softmax

$$\hat{y} * (1 - \hat{y})$$

3 COURSE - 3

WEEK 1

INTRODUCTION TO ML STRATEGY

① WHY ML Strategy

Things we can do to improve ~~problem~~ accuracy of model:

- ① Collect more data
- ② Collect more training, test data
- ③ Train longer with Gradient descent
- ④ Try bigger, smaller network
- ⑤ L2 regularization
- ⑥ change network architecture
- ⑦ dropout

② ORTHOGONALIZATION

* Some deep learning developers knows exactly what parameters to tune in order to try to achieve one effect. It is called orthogonalization.

* You'll have to fit training set well on loss function
 • If it's not achieved you could try bigger network,
 another optimization (like adam)

 • Fit dev set well on loss function.

 • If it's not achieved you could try regularization
 . bigger learning etc.

- iii) Fit test set well on cost function
 - If its not achieved you can try bigger dev. set
- iv) Performs well in real world.
 - If its not achieved you could try could try change dev. set change cost functions.

SETTING UP YOUR GOAL

① SINGLE NUMBER EVALUATION METRIC

- Its better and faster to set a single number evaluation metric for your project before you start it
- What are precision and recall
- o Suppose we run the classifier on 10 images which are 5 cats & 5 non-cat. The classifier identifies that there are 4 cats, but it identified 1 wrong cat
- o Then the confusion matrix will looks like this

	Predicted cat	Predicted non-cat
Actual cat	3	2
Actual non-cat	1	4

- o Precision : Percentage of true cats in the recognition result : $P = 3/(3+1)$
- o Recall : Percentage of true recognition cat of all cat predictions : $R = 3/(3+2)$
- o Accuracy : $(3+4)/10$

- Using precision/recall for evaluation is good in lot of cases, but separately they don't tell you which algorithm is better. But one thing we can do is combining combining the precision and recall called F1 Score.
- You can think of F1 Score as an average of precision and recall

$$F1 = 2/(C(1/P) + (1+R))$$

② SATISFYING AND OPTIMIZING

METRIC

Classifier	Precision	Recall	F1 Score
A	95%	90%	92.4%
B	98%	85%	91.0%

② SATISFYING AND OPTIMIZING

METRIC

Classifier	Accuracy	Running time
A	90%	80 ms
B	92%	95 ms
C	95%	1500 ms

- optimizing is the F1 Score
- satisfying is the running time of that particular accuracy. It has to be less than 100 ms

③ TRAIN / DEV / TEST DISTRIBUTION

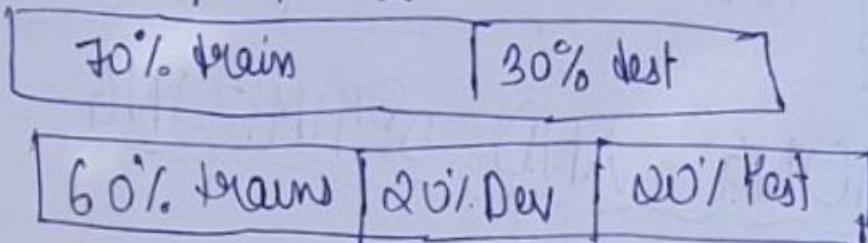
iii)

following are the guidelines for train/dev/test

- IV) ① ^{always} choose a dev set and test set from a same distribution
- ② and reflect data you expect to get in the future and consider important to do well on.
- ③ Setting up the dev set, as well as the validation metric is really defining what target you want to aim at.

④ SIZE OF THE DEV AND TEST SETS

→ Old way of splitting data



→ But in now modern machine learning era we are now used to working with much larger data set sizes. Suppose we 1,00,00,000 data example then we can use 98% to train our model, 1% for development and 1% for testing.

1% of 1,00,00,000 is 10000 which is quite high

5) WHEN TO CHANGE DEV/HESK SETS AND METRICS.

Example : (Cat classification)

Metric : classification error

Algo A : 3% error (But lot of poor images are treated as cat image here)

Algo B : 5% error

* we can easily say that our algo A is deviating by looking at error value but lot of poor images are also treated as cat images here

- * Thus, in this case, we ~~want~~ want and need to change our metric
 - ↳ we put less weight on images which are not poor
 $w[i] = 1$ if $x[i]$ is not poor
 - ↳ we'll put more weight on images which are poor
 $w[i] = \frac{100}{+}$ if $x[i]$ is poor
any big number
- * This is actually a example of an orthogonalization where you should take a Machine learning problem and break it into distinct steps:
 - i) Figure out how to define a metric that captures what you want to do - place that target
 - ii) Worry about how to actually do well on this metric - how to aim/shoot accurately at the target.

* If doing well on your metrics + dev/test set doesn't correspond to doing well or in your application, change the metrics or dev/test or both.

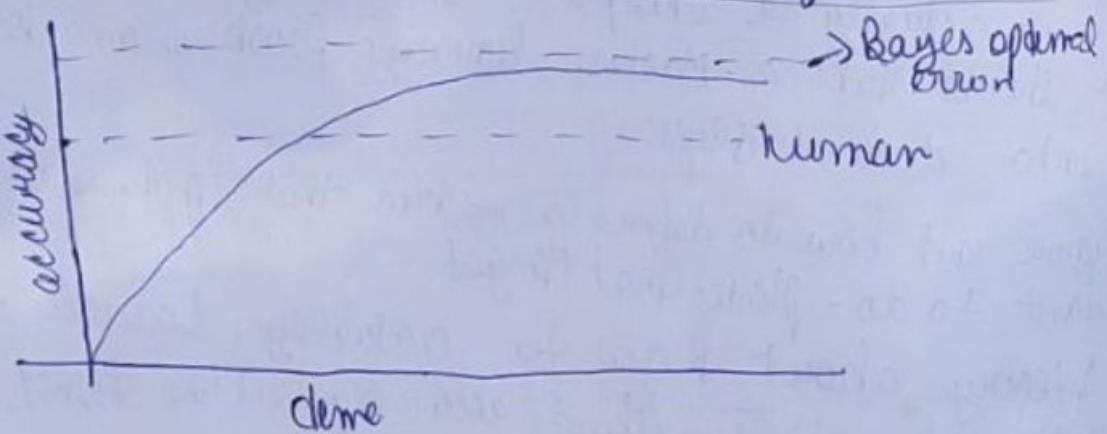
COMPARING TO HUMAN-LEVEL PERFORMANCE

① WHY HUMAN-LEVEL PERFORMANCE

Because of advances in deep learning, ML algorithms are suddenly working much better and so it has become much more feasible in a lot of application areas for machine learning algorithms to actually become competitive with human level performance.

It turns out that the workflow of designing and building a machine learning system is much more efficient when you're trying to do something that human can also do.

Comparing to human-level performance



Explanation of above diagram :-

- ① The more we fine tune ^{and}, the accuracy of an algorithm increases and it surpasses human level performance. ~~at after that~~ ~~it~~ we surpass human level performance there's not that much head room to still improve.
- ② A line will never meet Bayes optimal error because suppose we have blurry picture whether it is a cat or not even our algorithm can't tell whether it is cat or not.
- ③ Humans are quite good at lot of tasks so as long as ML is worse than humans, you can:
 - ① Get labeled of images by humans
 - ② Get insight from manual error analysis: Why did a person get it right?
 - ③ Better analysis of bias/variance.

2) AVOIDABLE BIAS

Suppose that a cat classification algorithm gives these results:

	I^{th}	II^{th}
Humans error	1%	7.5
① Training error	8%	8%
② Dev error	10%	10%

* maybe we has 7.5% of human error & 8% training error & dev error = 10%.
because of blurry images.

* Suppose we have 1% human error on I^{th} we can easily say that our model is not very good.
* If we have 7.5% (II^{th}) human error & 8% training error & dev error = 10%, we can say that our model is very good.

- * In the left example (IV. Human error) we have to focus on the bias
- * In the right example (7.5 Human error) then we have to focus on the variance.
- ↳ The human-level error is a proxy (estimate) for Bayes optimal error. Bayes optimal error is always less but human-level in most cases is not far from it.
- ✓ Also you can't do better than Bayes error unless you are overfitting
- ↳ Available bias = Training error - Human error
- * Variance = Dev error - Training error

(3) UNDERSTANDING HUMAN-LEVEL-PERFORMANCE

- * We have multiple human-level performance based on the human experience. Then you choose the human-level performance (proxy for Bayes error) that is more suitable for the system you're trying to build.
- * Human-level error (proxy for Bayes error)
 - Calculate available bias = training error - ^{human level error}
 - If available bias difference is ^{the} bigger, then bias problem
- * Training error
 - Calculate variance = dev error - training error
 - If variance difference is bigger then

You should use a strategy for variance reduction

- So having an estimate of human-level performance gives you an estimate of Bayes error and thus allows you to more quickly make decisions as to whether you should focus on trying to reduce a bias or try to reduce the variance of your algo.
- These techniques will tend to work well until you surpass human-level performance.

⑨ SURPASSING HUMAN-LEVEL PERFORMANCE

CE

Example :-

Team of humans	0.5%
One human	1%
Training error	0.6%
Dev error	0.8%

$$\text{ausable bias} = \text{difference b/w Team of Humans \&} \\ \text{training error} \\ = 0.1\%$$

$$\text{variance} = \text{difference b/w training error \& Dev errors} \\ = 0.2\%$$

In some problems, deep learning has surpassed human-level performance like:-

- Online advertising
- Product recommendation
- Image classification

- Humans are far better in natural perception tasks like computer vision & speech recognition
- But A.I. are doing better in some cases like - dangerous cancer, skin diseases etc

5) IMPROVING YOUR MODEL PERFORMANCE

Two fundamental assumptions of Supervised Learning
 IS →

- ① You can fit the training set pretty well with low avoidable bias
- ② The training set performance generalizes pretty well to dev/test set with variance is not too bad.

- * difference b/w human level & training set error is called avoidable bias
- * difference b/w the dev/test set and training set error is called variance

- o If avoidable bias is large you have these options:

- ① Train bigger model
- ② Train longer / better optimization algo (Momentum, RMSprop, adam)
- ③ find better NN Architecture / hyperparameters

or search.

- If variance is large you have these options:
 - Get more training data.
 - Regularization (L^2 , dropout, data augmentation)
 - Find better NN architecture / Hyperparameter search

WEEK - 2

ERROR ANALYSIS

① CARRYING OUT ERROR ANALYSIS

Error analysis refers to the process of manually examining mistakes that your algorithm is making.

Example:

- You discovered that some of the mislabeled data are dog pictures that look like cats. Should you try to make your cat classifier do better or dogs
- Error analysis approach:
 - Get 100 mislabeled dev set examples at random.
 - Count up how many are dogs.
 - If 5 of 100 are dogs then training your classifier to do better on dogs will decrease your error upto 95% which can be too little.
 - If 50% of 100 are dogs then you could decrease your error up to 5% which is reasonable and you should work on that.
- Based on last example, error analysis helps you to analyse the errors before taking an action

that could take lot & come with no need.

- Sometimes, you can evaluate multiple cross analysis ideas in parallel & choose the best idea. Create a spreadsheet to do that and decide eg:

	Image	Dog	Great cats	Blurry	Auto filters	Commands
1		✓				Piobal
2		✓		✓	✓	Rainy day at Zoo
3						
4			✓			
	% total	8%	43%	61%	12%	

- In the last example you will decide to work on great cats or blurry images to improve your performance
- This quick counting procedure, which you can often do in at most, small numbers of hours can really help you make much better prioritization decisions, and understand how promising different approaches are to work on.

② CLEANING UP INCORRECTLY LABELED DATA

If your learning algorithm is not yet at the performance of a human. Then manually examine mistakes that your algorithm is making can give you insights into what to do next. This process is called error analysis

* If you are working on a brand new ML application then you should build your first system quickly and then iterate.

- ¶
 - ① Set up dev/test dev and metrics
 - ② ~~Build~~ Build initial system quickly
 - ③ Use bias/variance analysis & error analysis to prioritize next steps.

MISMATCHED TRAINING AND TEST SET

① TRAINING AND TESTING ON DIFFERENT DISTRIBUTIONS

Suppose we want to build a mobile application in which we want in which a application will show to tell whether the photo clicked by user is a cat or not.

One strategy we can do is
Following are the two strategies-

- ① Option one (not recommended)
In this we train our model using the pictures we find of cat on internet and take picture clicked by user and then we shuffle and then we train our model but in this strategy ~~K~~ ~~train~~ accuracy will be low and even if our model is not good with pictures pictures clicked by user it will show high accuracy.

② Option two (Recommended)

Train your model using pictures of cats from internet, and then we can use picture of cats clicked by users on test sets. This will give us the real accuracy of the model.

③ BIAS AND VARIANCE WITH MISMATCHED DATA DISTRIBUTIONS

The way you analyze bias and variance changes when your training set comes from different distribution than your dev and test sets.

- The cat classifier example. Suppose you worked in the example and reached this
 - Human error - 0%.
 - Train error - 1%.
 - Dev error = 10%

In this example, you'll think that this is a variance problem, but because the distribution aren't the same you can't tell for sure, because it could be that train set was easy to train on, but the dev set was more difficult.

Note variance difference is 9%, which is big but it might be good.

- So, the problem with this analysis is that you went from one training error to dev error due things changed at a time.

- One is that the algorithm sees data in the training set but not in the dev. set.
- Two, the distribution of data in the dev set is different.

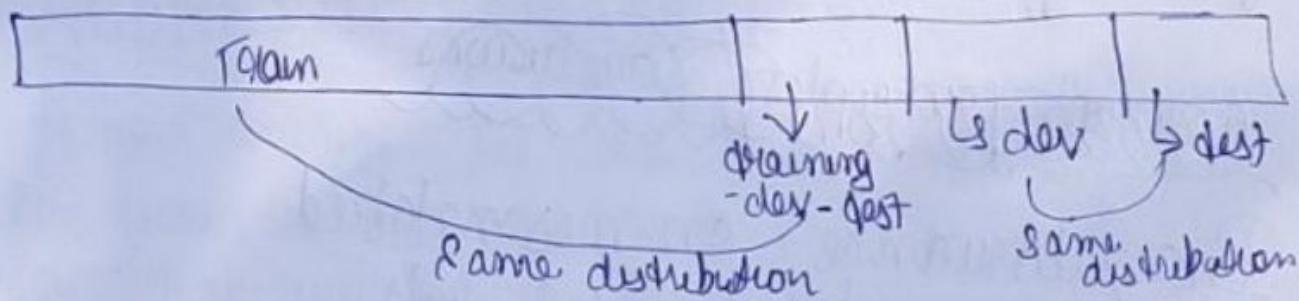
and because you changed two things at the same time, it's difficult to know of this 9% increase in error, how much of it is because the algorithm didn't see the data in the dev set and how much of it, is because the dev set data is just different.

In order to tease out those two effects it will be useful to define a new piece of data which we'll call the training-dev dev set (now dubbed 3) data, which we come out that should have the same distribution as training set.)

Training-dev dev set: same distribution as training set, but not used for training

example

Now we're going to do is randomly shuffle the training sets and then come out just a piece of the training set to be the training-dev dev set



To carry out error analysis, what you should do is now look at the error of your classifier on the

training set, on the training-dev set, as well as on the dev set.

Example:

Now like this:

Human error : 0%.

Train error : 1%.

Train-dev error : 9%.

Dev error : 10%.

Now we can see that this a high variance problem

another example,

Human error - 0%.

Train error - 1%.

Train-dev error - 1-5%.

Dev error - 10%.

In this case we have something called data mismatch problem because your learning algorithm was not trained explicitly on data from training-dev or dev but these two data sets come from different distribution.

~~Generalization vs Conclusions~~

Bias/variance on mismatched training and dev/test sets

Human level	4%] - avoidable bias
Training dev error	7%] variance
Training dev dev error	10%] data mismatch
Test error	12%	degree of overfitting degree to dev set.

- * If the difference is big then its avoidable bias problem other you should use a strategy for high bias
- * If the difference is big then its high variance problem then you should use a strategy for dealing it. (variance)
- * If difference is much bigger than train-dev error its data o mismatch problem.
- * degree of overfitting to dev set = test error - dev error
 - As the difference is big (positive) then & maybe you need to find a bigger dev set (dev set and test set come from the same distribution)

③ ADDRESSING DATA MISMATCH

If your training set comes from a different distribution than your dev and test set, and if error analysis shows you that you have a data mismatch problem.
what to do?

There is no completely systematic solution to this, but there some things you could try.

Addressing data mismatch

- Carry out manual color analysis to try to understand difference b/w training & dev/test sets
- Make training data more similar; or collect more data similar to dev/test sets

If your goal is to make the training data more similar to your dev set one of the techniques you can use artificial data synthesis that can help you make more training data

↳ combine normal audio with car noise to get audio with car noise example

↳ But we do need to be cautious and bear in mind whether or not you might be accidentally generating data from a tiny subset of the space of all possible examples because your UN might overfit these generated data

LEARNING FROM MULTIPLE TASKS

① Transfer learning

Transfer learning refers to using the NN knowledge for another application.

When do we use transfer learning?

- ① Task A & B have same input x .
- ② A lot more data for task A than task B
- ③ Low-level features from task A could be helpful for task B

Example 1:

The following network is trained for cat recognition, but we want to adapt it for radiology diagnosis. The neural network will learn the structure & the nature of images. This initial phase of training on image recognition called pre-training, since it will pre-initialize the weights of NN. Updating all the weight afterwards is called fine-tuning.

Example 2:

Let's say that you've trained a speech recognition system so now x is input of audio. If y is some text transcript. So, you've trained an speech recognition system to output your transcript. And let's say that you now want to build a wake words or trigger words detection system such as "Hey Alexa, Ok Google", to wake up device.

~~② Multi-task learning~~

② MULTI-TASK LEARNING

↳ In multi-task learning, you start off simultaneously trying to have one neural network do several things at the same time and then each of these tasks helps hopefully all of the other tasks.

~~↳~~ Example:

Let's say you're building an autonomous vehicle, building a self-driving car. Then your self-driving car would need to detect several different things such as pedestrians, detect other cars, detect stop signs. Then Y shape will be (Y, m).

We can also train 4 NN separately. But if some of the earlier features in neural network can be shared b/w those different objects, then you find that NN. do do 4 things

$$Y = \begin{bmatrix} 1 & ? & 1 & \dots \\ 0 & 0 & 1 & \dots \\ ? & 1 & ? & \dots \end{bmatrix}$$

Why ~~not~~ multi-tasking?

- ↳ Training on a set of tasks that could benefit from having shared lower-level features
- ↳ Usually, amount of data you have for each task is quite similar
- ↳ Can train a big enough network so do well on all the tasks.

If you train a big enough NN, the performance of the multi-task learning compared to splitting the task is better. Today transfer learning is used more often than multi-task learning.

TRANSITION FROM MULTITASK LEARNING END-TO-END DEEP LEARNING

Q What is end-to-end deep learning?

Some systems have multiple stages to implement. An end-to-end deep learning system implements all these stages with a single NN.

Example :-

- Speech recognition system

audio → features → phonemes → words → transcript # non-end-to-end system

audio → transcript # end-to-end learning system

- End-to-end learning gives data more freedom; it might not use phonemes when training.

Q To build the end-to-end deep learning system that

works well, we need a big dataset (more data than in non-end-to-end system).

If we have a small dataset the ordinary implementation could work just fine.

Whether to use end-to-end deep learning

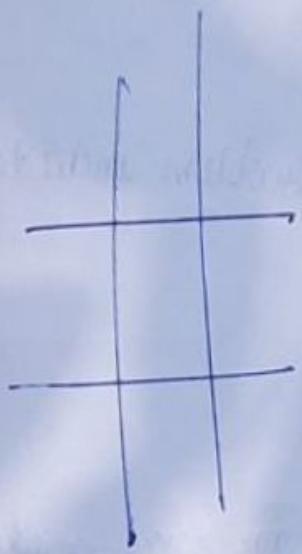
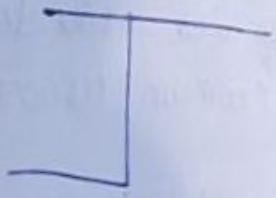
- Pros of end-to-end deep learning
- ~~drawbacks~~ By having a pure ML approach, your NN learning input from X to Y may be more able to capture whatever statistics are in the data, rather than being forced to reflect human ~~preconceived~~ preconceptions.
- less hand-designing of components needed.

Cons of end-to-end learning

- May need a large amount of data.
- Excludes potentially useful hand-design components

Applying end-to-end deep learning:

- Key question? Do you have sufficient data to learn a function of the complexity needed to map X to Y ?
- Use SML/DL to learn some individual components
- When applying supervised learning you should carefully choose what types of X to Y mapping task you want to learn depending on what you can get data for:



#4

CONVOLUTIONAL NEURAL NETWORKS

WEEK 1

1 COMPUTER VISION

What is computer vision?

Computer vision is a field of computer science that works on enabling computer to see, identify and process in the same way that human vision does.

Some of the applications of Computer Vision that are using deep learning are:

- Self driving cars.
- Face recognition.

Examples of a computer vision problem include

- Image classification
- Object detection
 - Detect object & localize them
- Neural style transfer

▫ Changes the style on an image using another image (filters), adding ~~animate~~ (making) ~~animate~~ ~~animate~~ to make the

Suppose we have a picture of cat, that $64 \times 64 \times 3$ ($64 \times 64 \times 3$)
 Suppose we have a cat picture of $64 \times 64 \times 3$ ($64 \times 64 \times 3$) that means we have 12288 input features.
 But what if we have $1000 \times 1000 \times 3$ image? ($1000^2 = 1 \text{ Mega Pixel}$)
 If we have 3 million input features than this means that x will be 3 million dimensional.

To do that you

One of the solution is to build this using convolution layers instead of the fully connected layer.

EDGE DETECTION EXAMPLE

In this we'll see how convolution operation works.

Example:

My photo of object

Given a picture like that for a computer do figure out what are the objects in this picture. The first thing you'll do is maybe detect vertical edges in this image and we also want to find horizontal edges. So how we detect edges in an image?

Here is an 6×6 grayscale image

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

An order to detect edges or lets say vertical edges in this image we will construct a 3×3 matrix known as filter

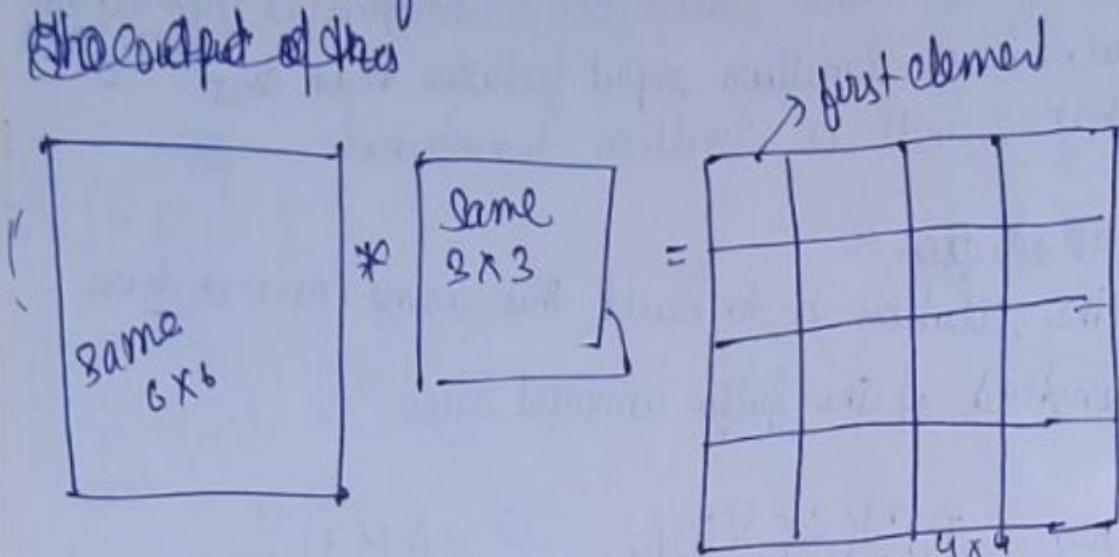
*

1	0	-1
1	0	-1
1	0	-1

Sometimes filter is also called Kernel

and what you are going to do is take the 6 by 6 image and convolve it and the convolution operation is denoted by asterisk (*) (this is not multiplication). This asterisk is referred to convolution.

The output of this



The output of the convolution operator will be a 4×4 matrix.

The way to compute this 4×4 output is as follows :

* The first element of 4×4 matrix is calculated as what you are going to do is take the 3×3 filter and place it on top of the 3×3 region of your original input image. All this

3	0	0	-1			
1	5	0	8	-1		
2	7	0	2	-1		

$$3x_1 + 1x_1 + 2x_1 + 0x_0 + \\ \dots = -5$$

$6 \times 6 \rightarrow$ same matrix

and to compute 2^{nd} element we do after the same calculation but now we take one step right like this :-

0	1	2	-1	7	4
1	8	9	-1	3	1
5	1	2	5	-1	1
7	1	2	5	-1	3
1	1	1	1	1	1

*
↓
conv

1	0	1
1	0	-1
1	0	-1

$$= \begin{bmatrix} -5 & -4 & - \\ & & \end{bmatrix}$$

This is how we calculate all the value and this turns out to be a vertical edge detector

We can do this ~~tensorflow~~ by using this function
tf.nn.conv2d

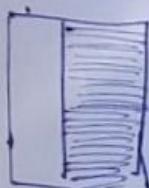
Let's look at another example :-
and you'll get

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

$$\ast \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} =$$

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

If we plot this image it might look like this



if filter can
be written as
follows

30	30	0
30	30	0
30	30	0



In our dimensions row seems little bit wrong that the detected image edge seems really thick, that's only because we are working with very small images in this example. Will generally use say 1,000 by 1,000 px image rather than 6x6 images when you find that this does a pretty good job.

② MORE DETECTION EXAMPLE

Vertical and Horizontal Edge detection

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

Vertical

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

Horizontal

example for horizontal edge detection

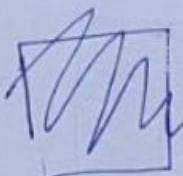
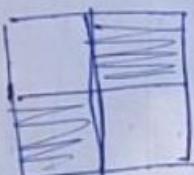
$$\begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \end{bmatrix}$$

*

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

=

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 30 & 0 & -10 & 30 \\ 30 & 10 & -10 & -30 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$



So, in summary, different filters allow you to find vertical and horizontal edges. It turns out that the three by three vertical edge detection filter value used is just one possible choice.

In computer vision literature there was a fair amount of debate about what is the best set of numbers to use. So here's some of it.

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 2 & 0 & -2 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$

Sobel
filter

$$\begin{array}{|c|c|c|} \hline 3 & 0 & -3 \\ \hline 10 & 0 & -10 \\ \hline 3 & 0 & -3 \\ \hline \end{array}$$

Scharr filter

And this is just for vertical edge detection if we flip this it will become horizontal edge detector.
 $\Rightarrow 90^\circ$

When you really want to detect edges in some complicated image, maybe you don't need to have computer vision researchers handpick these nine numbers. Maybe you can just learn them and treat the nine numbers of this matrix as parameters, which you can then learn using back propagation and the goal is to learn nine parameters so that when you take the image the 8×8 image and convolve it with your three 3×3 filter then this will give you a good edge detector and by just letting all of these numbers be parameters and learning them.

automatically from data, we find that NN can actually learn low level features.

③ PADDING

If we have 6×6 image and 3×3 filter then
the dimension of the output will be
 $n-f+1 \times n-f+1$
 $6-3+1 \times 6-3+1$
 4×4 (output matrix)

Note our image shrinks after applying filter and also
the pixel at the corner or the edge is used only
in one of the output (which is 1st element for apply
filter for the first time) whereas if the filter
is applied many times
pixels on the corners or on the edges are
used much less than pixels in the middle. So here
we are throwing away lot of information.

To solve this problem we use padding.
padding refers to adding padding to the image
with additional border.

0	0	1	0	1	0	0
0	5	0	7	7	0	.
0	6	7	3	9		
0	8	2	1	2	0	
0	7	1	3	1	0	
0	9	7	7	1	0	
0	0	0	0	0	0	

→ This will change the pixel
dimension of input
~~4x5~~ ~~5~~ 6
to
6x7 after padding

So how to calculate the dimension of output after padding?

Here is formula for this.

$$n+2p-f+1 \times n+2p-f+1$$

{ p = padding

n = input

f = filter

So after padding edges / corner pixel are used many times for calculation of output.

We can all do padding of 2 pixel. This will add 2 borders to our input pixel.

In terms of how much to pad, it turns out there are two common cases that are called Valid and Same convolutions.

Valid convolution \rightarrow No padding

input will be $n \times r$ filter will be $F \times F$

and output dimension will be $n-f+1 \times n-f+1$

Same convolution \rightarrow Pad to that output size

is same as the input size.

when we pad with p pixel

$$n+2p-f+1 \times n+2p-f+1$$

$n+2p-f+1 = n$ (so the output size is same as input size)

$$= p = \frac{f-1}{2}$$

f is usually odd and you'll rarely see even numbered filters.

We have 2 reasons for that :-

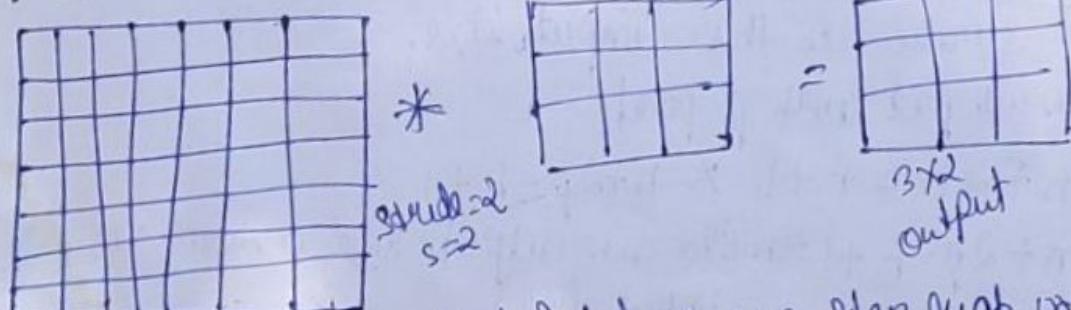
- ① If f was even then you need some asymmetric padding, so only if f is odd that this type of same convolution gives a natural padding region. (Pad more on the left pad more less on right or opposite called asymmetric padding)
- ② When you have an odd dimension filter, such as 3×3 then it has a central position and sometimes in computer vision its nice to have a distinguisher (central pixel)

Maybe none of this is a great reason to have 3×3 filters, 3×3 are very common

It will give fine performance even if we use 4×4 filter

⑨ STRIDED CONVOLUTION

We want to convolve this 7×7 image by 3×3 filter with stride of 2



In this instead of taking one step right we do two steps

If we stride at 2 ($S=2$) then what is the dimension of the output? Here the formula

$$\frac{n+2p-f}{S} + 1 \quad \times \quad \frac{n+2p-f}{S} + 1$$
$$\frac{7+0-3}{2} + 1 \quad \times \quad -$$
$$= \underbrace{3}_{\substack{\uparrow \\ \text{if this value is not int in that case will}}} \times 3$$

If this value is not int in that case will round this down

$$\lfloor z \rfloor = \text{floor}(z)$$

\rightarrow floor symbol

Your 3×3 filter must be entirely within your image (in image + padding)

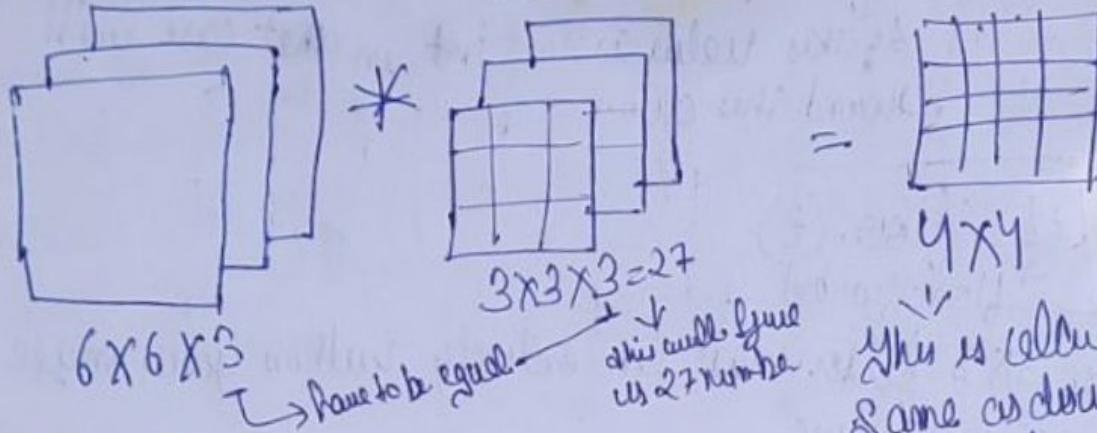
Technical note on cross-correlation vs convolution

In math textbooks the conv operation is flipping the filter before using it. what we were doing is called cross-correlation operation but the state of art of deep learning is using this as conv operation.

Technically, what we're actually doing the operation we've been using for the last few videos is sometimes cross-correlation instead of convolution. But in deep learning literature by convention we just call this a convolutional operation.

⑤ CONVOLUTIONS OVER VOLUME

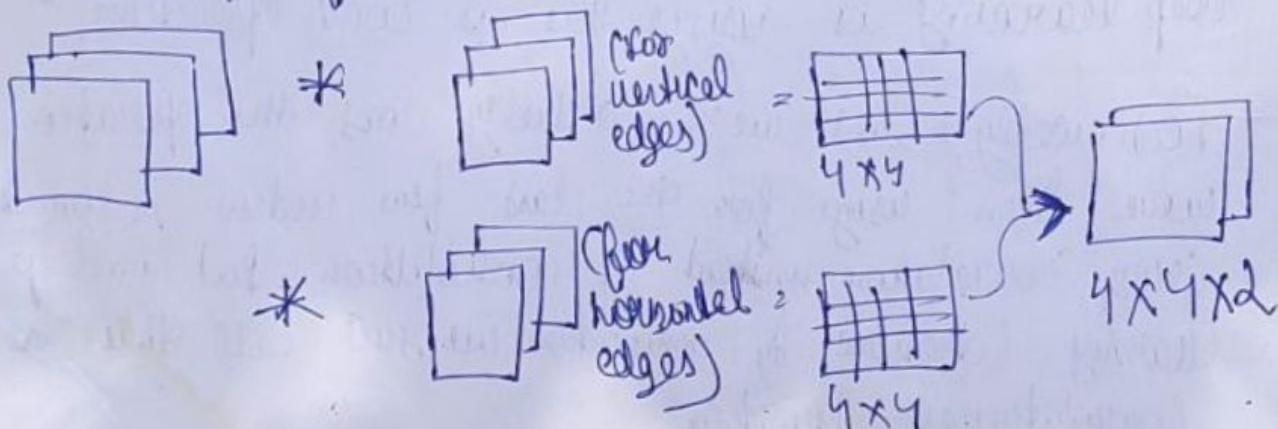
Let's say you want to detect features, not just gray scale image but in RGB image like $(6 \times 6 \times 3)$ then we'll have feature image of $(3 \times 3 \times 3)$



Now we have learned to convolve on volumes.

This is calculated same as discussed but for every layer. Then add up all the numbers, this will give you your number.

What if we don't just want to detect vertical edges? What if we wanted to detect vertical edges and horizontal edges in short what if you want to use multiple filters at the same time?



(5)

Summary:

$$\begin{aligned} \text{Input} &= n \times h \times w \xrightarrow{\text{number of channels}} n_c \\ &= 6 \times 6 \times 3 \end{aligned}$$

$$\begin{aligned} \text{filter} &\xrightarrow{\text{number of channels}} n_c \\ \text{Filters} &\xrightarrow{f \times f \times n_c} \\ &\quad 3 \times 3 \times 3 \end{aligned}$$

= output

$$n-f+1 \times n-f+1 \times n_c$$

$$4 \times 4 \times 2$$

no. of filters
↑

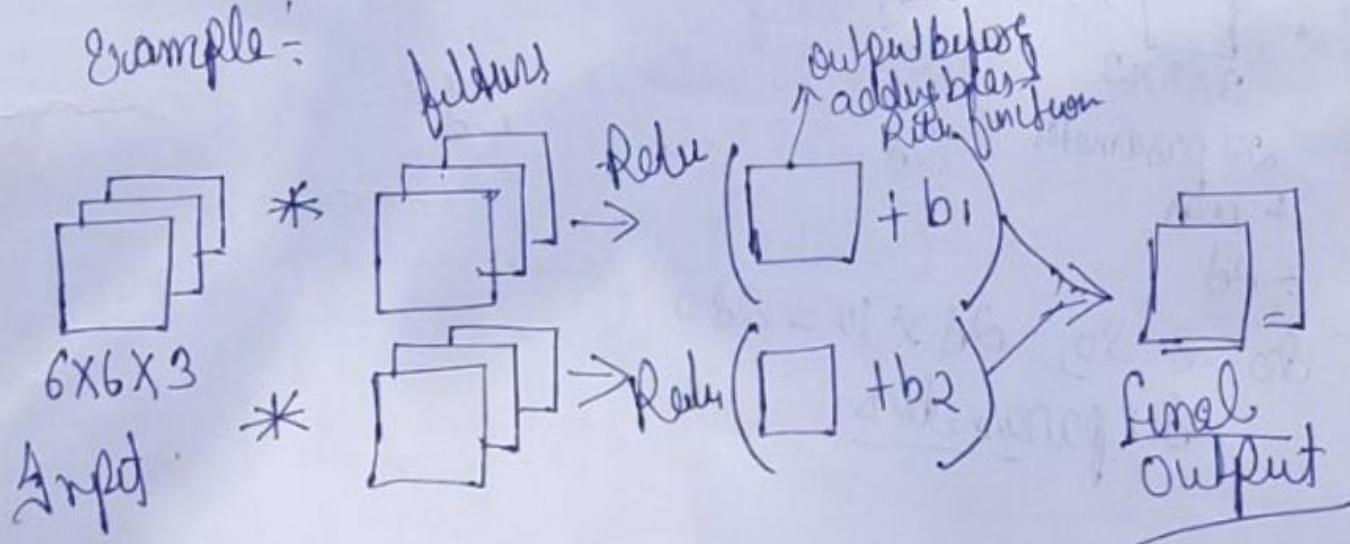
| # no. of input should be same
as no. of filters

If we use stride more
than one (1) then this formula
will change

- * The output have number of channels equal to the number of filters you are detecting.
- * The last dimension in input layer in M.L. literature people often call this the depth of this 3D volume

⑥ ONE LAYER OF A CONVOLUTIONAL NETWORK

Example:



This is one layer of convolutional network

So to map this back to one layer of four propagation
in the standard N.N. in a non-convolution network
remember.

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

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

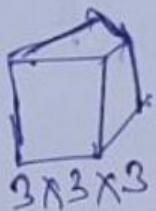
- * filters are similar to $w^{(1)}$
- * $a^{(1)}$ are similar to input
- * $b^{(1)}$ are bias
- * $a^{(1)}$ is output

Number of parameters in one layer

Q) If you have 10 filters that are $3 \times 3 \times 3$ in one
layer of N.N. how many parameters does that layer

have

Sol]



27 parameters

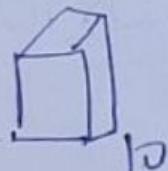
+ bias

= 28



28

- - -



280

So so so, $28 \times 10 = 280$

280 parameters

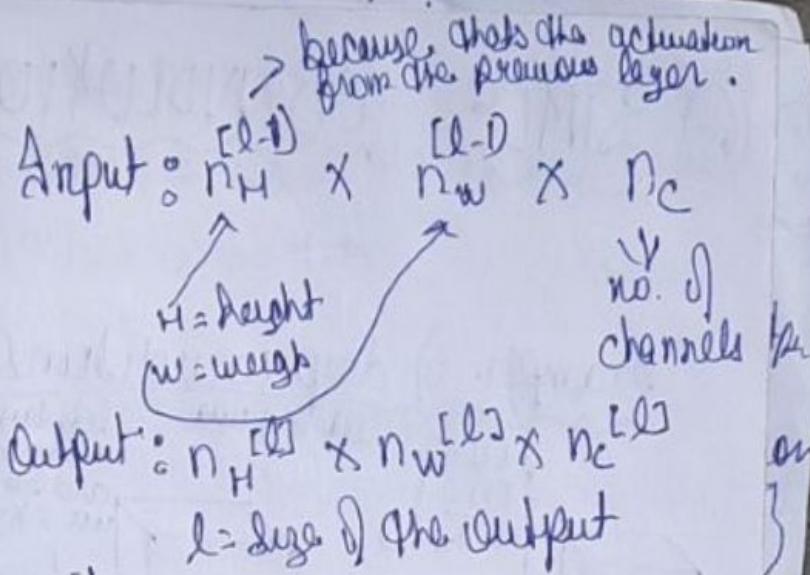
Summary of notation

f^2 = filter size

p^2 = padding

s^2 = stride

$n_c^{[l]}$ = number of filters



Each filter is: $F^{[l]} \times f^{[l]} \times n_c^{[l]}$

Activations: $a^{[l]} \rightarrow n_H^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

and when you are using a vectorized implementation or batch gradient then your actually output a is -

$A^{[l]} \rightarrow m \times n_H^{[l]} \times n_w^{[l]} + n_c^{[l]}$
↓
m samples

Weights: $F^{[l]} \times f^{[l]} \times n_c^{[l-1]}$ (What's the dimension of one filter)

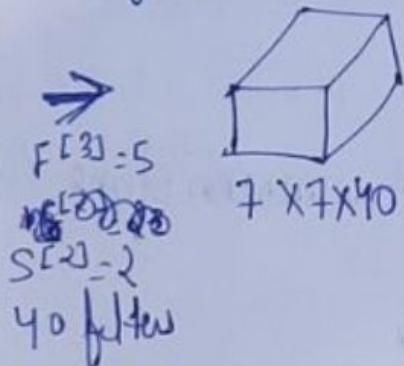
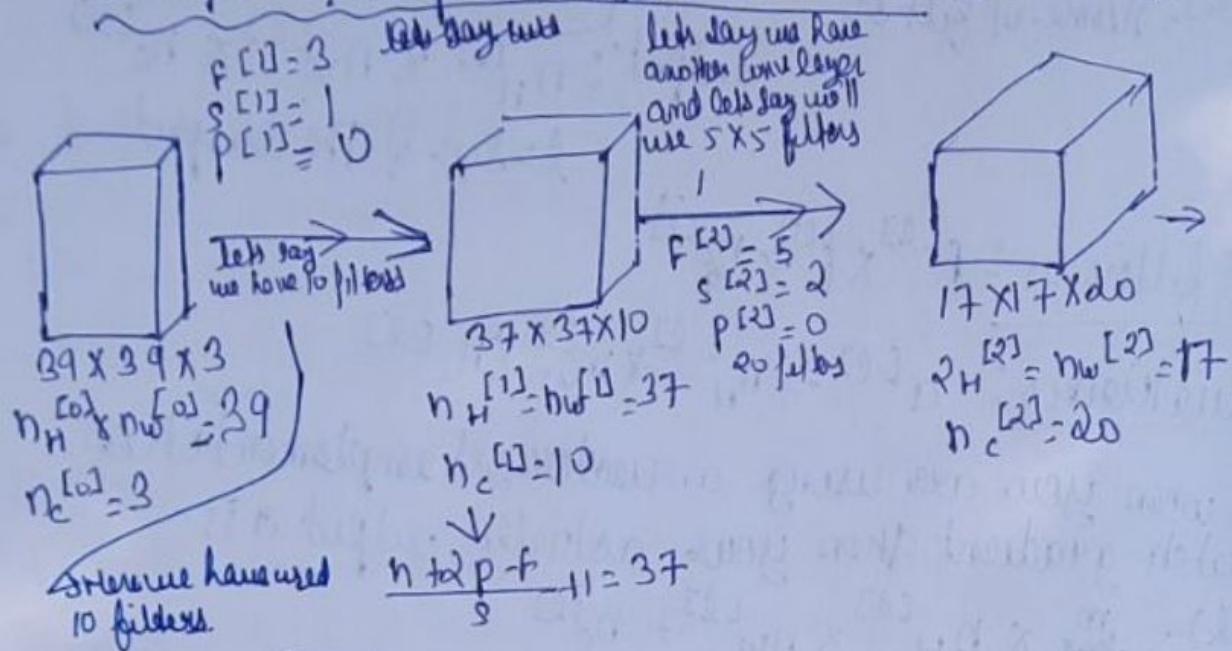
$n_c^{[l]}$ are the number of filters

& $F^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$

Bias: $n_c^{[l]} \rightarrow$ (Real number)

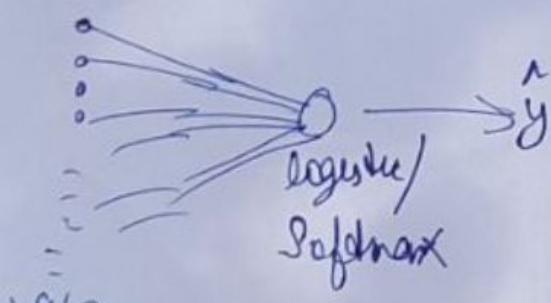
⑥ SIMPLE CONVOLUTIONAL NETWORK

example of deep convolutional network (ConvNet)



7 times 7 times 40 is 1960.

So what we can do is take this volume and flatten it.



1960

\hookrightarrow It is a very long vector

Note: the size of the images are tend to go down 39×39 to $3 \times 37 \times 37$ to 17×17 to ~~$10 \times 7 \times 7$~~ while the channels will generally increase 3 to 10 to 20 to 40

Types of layer in a convolutional network:

- convolution (Conv)
- Pooling (Pool)
- Fully connected

You can ^{also} create a convNet only using convolution

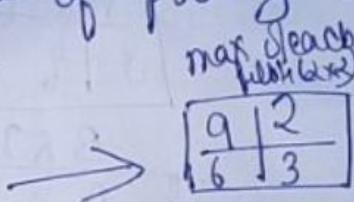
⑦ POOLING LAYERS

other than convolutional layers, ConvNets often also use pooling layers to reduce the size of the representation. ConvNets often use pooling layers to reduce the size of the representation, to speed the computation.

example:

Suppose you have a 4×4 input and you want to apply a type of pooling called max pooling.

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2



So, this is as if you apply a filter size of two
 $P=2$] hyperparameters
 $S=2$ of max pooling

So here's the intuition behind what max pooling is doing.
 If you think of this 4×4 region as some set of features,
 the activations in some layer of NN. Then a large
 number means that it may have detected a particular feature.
 So, what the max operation does is a lots of
 features detected anywhere and one of these quadrants,
 it then remains preserved in the output of max pooling.
 If those features are detected anywhere in this filter
 then keep a high number. But if this feature is not detected, keep low number. So that's the intuition
 behind max pooling.

One interesting property of max pooling is
 that it has a set of parameters (S and F) but it
 has no parameters to learn. Gradient descent
 will not change this parameter.

Pooling layer : Max Pooling

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

\rightarrow

9	9	5
9	9	5
8	6	9

3×3

$f = 3, S = 1$

At ii , called Samed in 2D input ~~and~~.

Feeding time : Average feeding



→



Averaging of feeding:

Hypothetical case:

1: After 3 hr. = 400000

2: 3 weeks

will be average feeding

case 1: generally average feeding & feeding together
will do it one together case will have the
same value

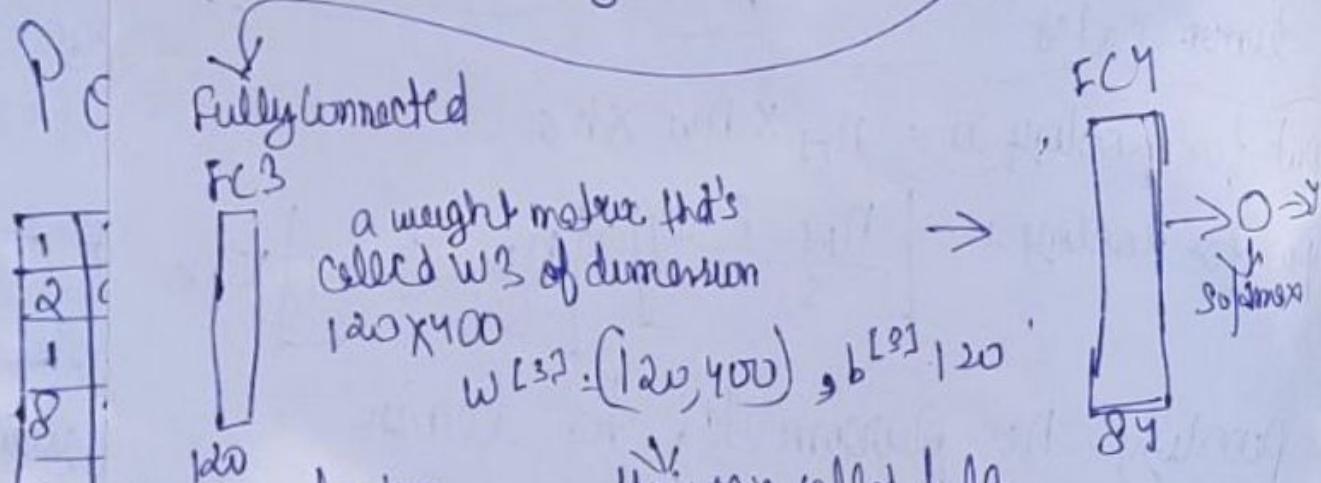
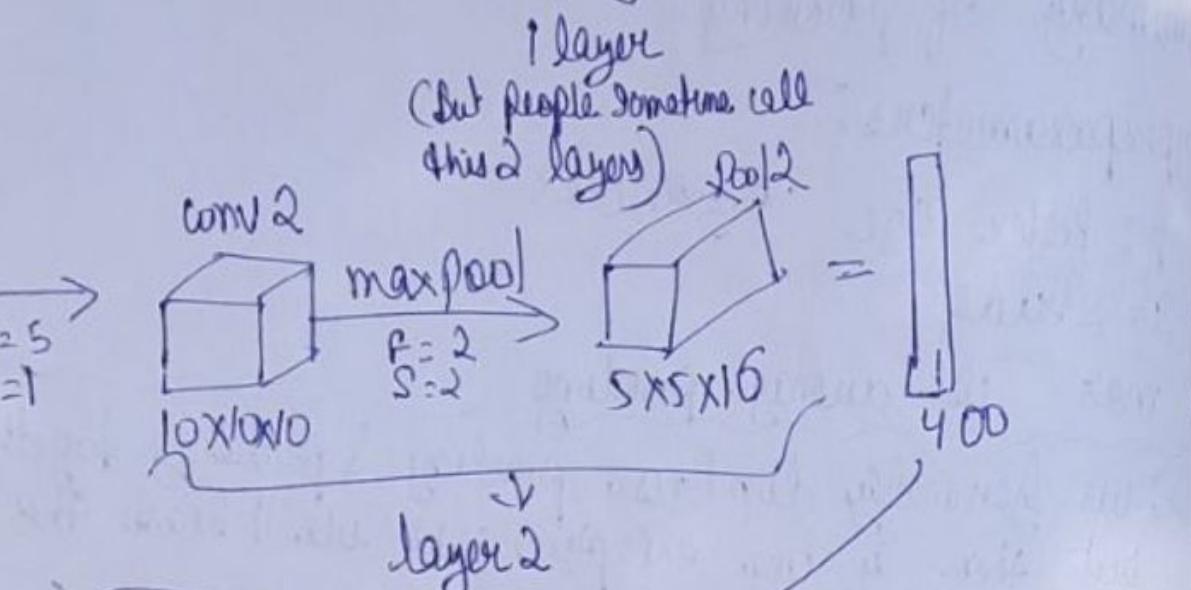
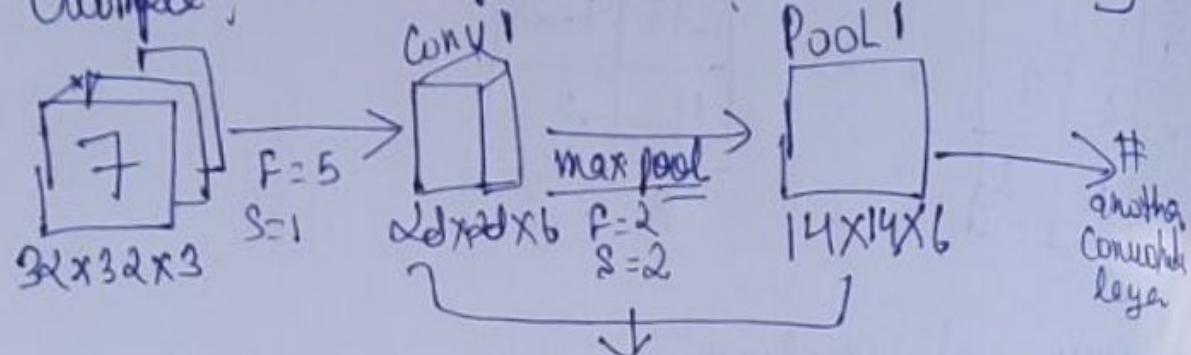
avg for feeding = $\frac{P_{11} + P_{12} + P_{13}}{3}$ %

avg for feeding = $\left[\frac{P_{21} - E}{3} + 1 \right] \times \left[\frac{P_{22} - E}{3} + 1 \right] \times \%.$

13. feeding no permission to start

⑧ CNN EXAMPLE

Example: (This example is something like LeNet-5)



Because we have 400 units densely connected to 120 units.

At single NN layer

This is called fully connected because each of the 400 units here is connected

Some statistics about the last example:

	Activation shape	Activation size	# Parameters
Input	(32, 32, 3)	3072	0
Conv (f=5, s=1)	(28, 28, 8)	6272	208
POOL 1	(14, 14, 8)	1568	0
Conv2 (f=5, s=1)	(10, 10, 16)	1600	416
POOL 2	(5, 5, 16)	400	0
FC3	(130, 1)	130	4800
FC4	(84, 1)	84	10,080
Softmax	(13)	10	841

⑨ WHY CONVOLUTIONS?

Advantages

① Parameter sharing

- A feature detector (such as vertical edge detector) that's useful in one part of the images is probably useful in another part of the images.

② Sparsity of connections

- In each layer, each output value depends only on a small number of inputs which makes it translation invariant

Paddings if all together:

WEEK 2

① WHY LOOK AT CASE STUDIES?

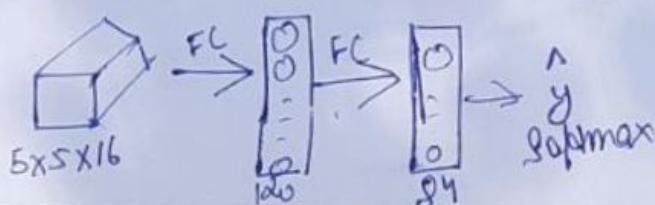
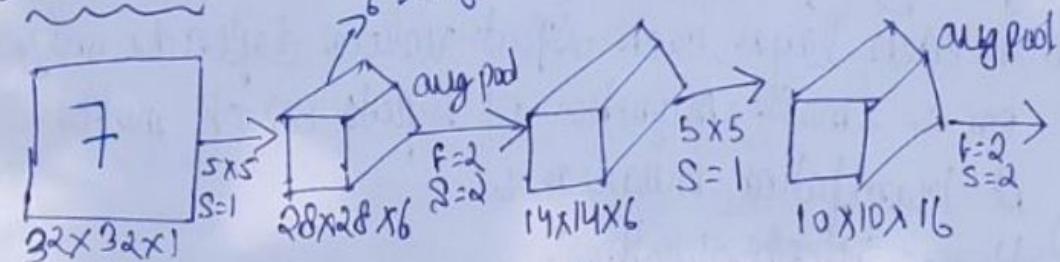
- Some NN architecture that works well in some tasks can also work well in other tasks.
- Here are some classical CNN network:
 - LeNet-5
 - AlexNet
 - VGG
- The best CNN architecture that won last ImageNet competition is called ResNet & it has 152 layers

② CLASSIC NETWORKS

Another some of the classic neural network are:

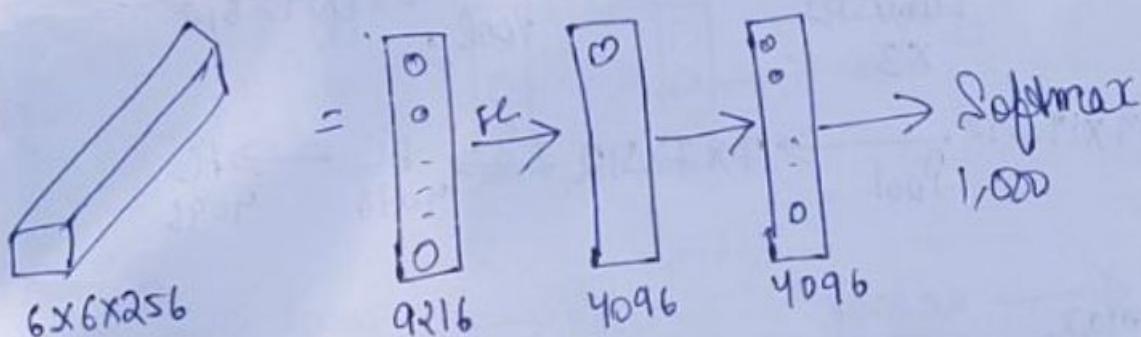
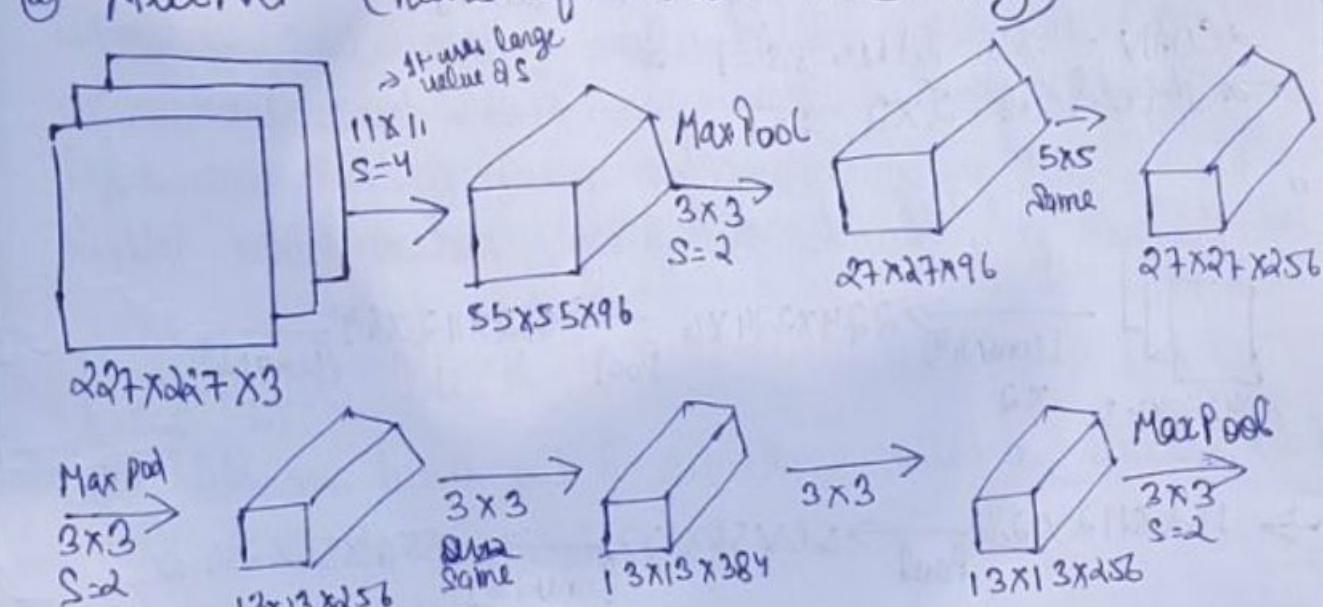
- ① LeNet-5
- ② AlexNet
- ③ VGG

① LeNet-5 *as it is trained on grayscale image*



- This model was published in 1998. The last layer wasn't using softmax back then
- It has 60K parameters
- Dimensions of the images ↓ as no. of channels ↑
- Conv → Pool → Conv → Pool → FC → FC → Softmax this type of arrangement is quite common
- The activation function used in the paper was Sigmoid & Tanh

② AlexNet (named after Alex Krizhevsky)

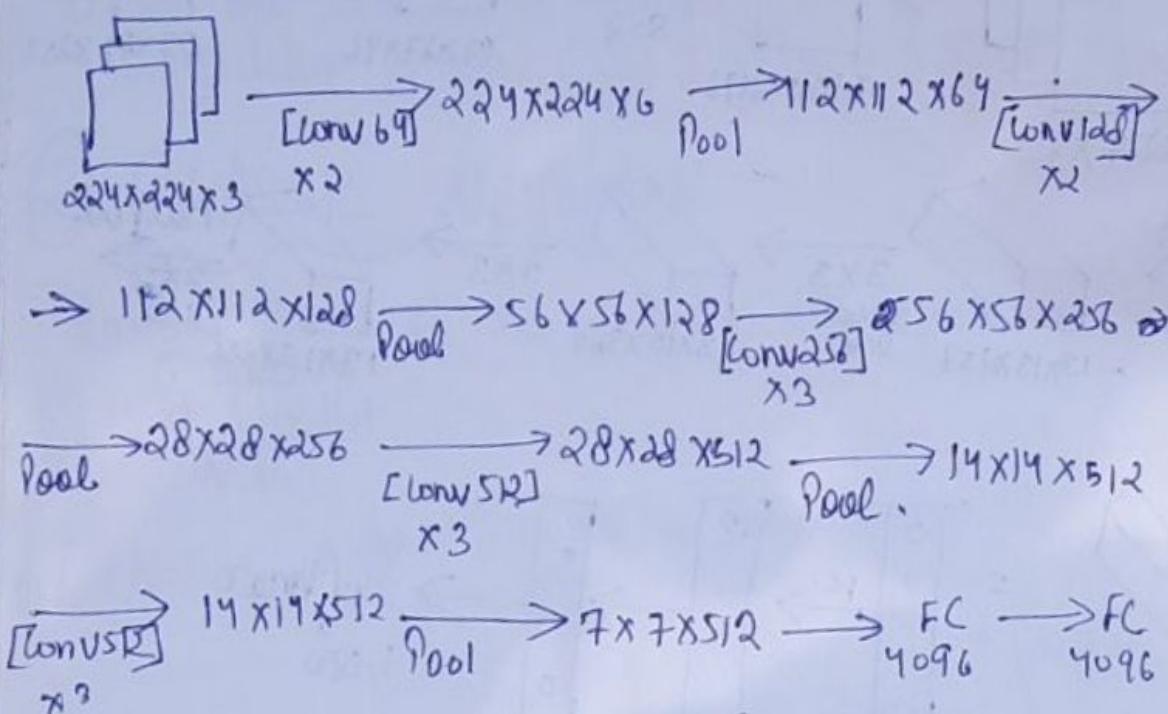


- Similar to LeNet-5 but has 60 million parameters compared to 60K parameters of LeNet-5
- Has ReLU^{activation} function

- Multiple GPUs were used because the GPUs were not so fast back then

⑨ VGG-16

- Instead of having so many hyperparameters^{it} use a much simpler network where you focus on just having CON-layers that are just 3×3 filters with a stride of one and always use same padding.
- It focuses on having only these blocks:
 - CONV = 3×3 filter, $s=1$
 - Max-Pool = 2×2 , $s=2$
-



\rightarrow Softmax
1000

- This network is large even by modern standards. It has around 138 million parameters.

- At has a total memory of 96MB per image for only forward propagation
- Number of filters increase from 64 to 128 to 256 to 512. 512 was made twice

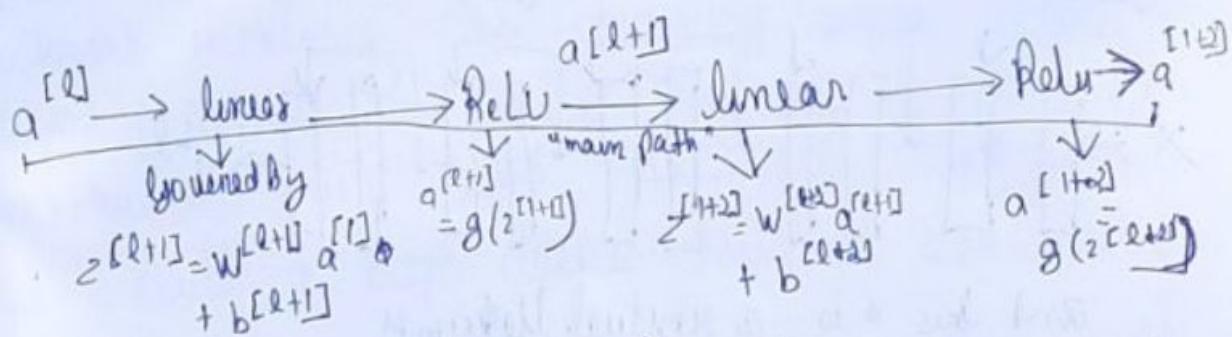
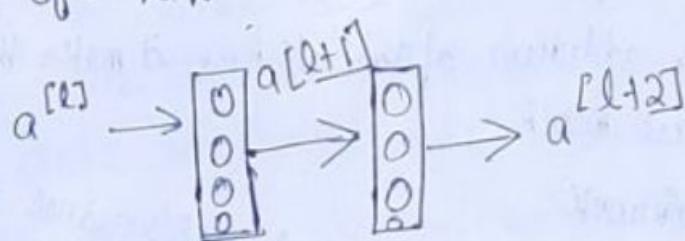
③ RESNETS

Very, very deep NN are difficult to train because of vanishing and exploding gradient. We'll learn about skip connections which allows you to take the activation from one layer and suddenly feed it to another layer. You'll build ResNet which enables you to train very, very deep networks.

RESIDUAL BLOCK

ResNets are built out of something called a residual block. Let's first describe what that is -

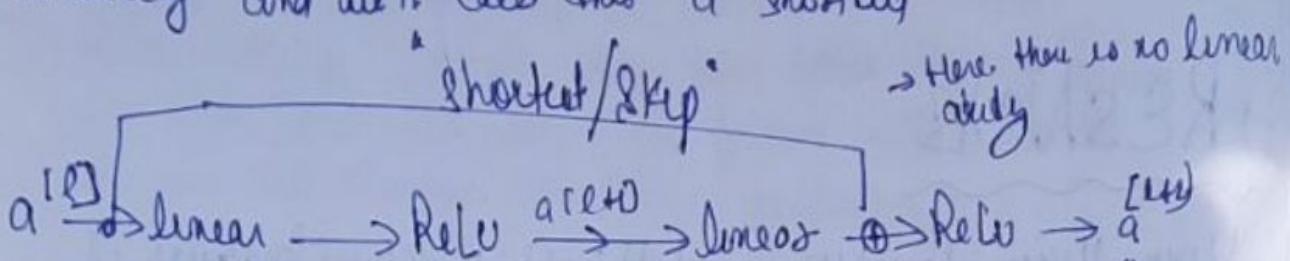
of layers of a NN



This is called main path, and we are going to change this

we're going to going make a change to this.

We're going to take $a^{[l]}$ and just forward it, and just add $a^{[l]}$ before applying to non-linearity, the ReLU non-linearity and we'll call this a shortcut



So, rather than needing to follow the main path the information from $a^{[l]}$ can now follow a shortcut to go much deeper into the N.N. also this means that last equation goes away and we instead have that the output $a^{[l+2]}$ is the ReLU non-linearity g applied to $z^{[l+2]}$ as before but now plus $a^{[l]}$

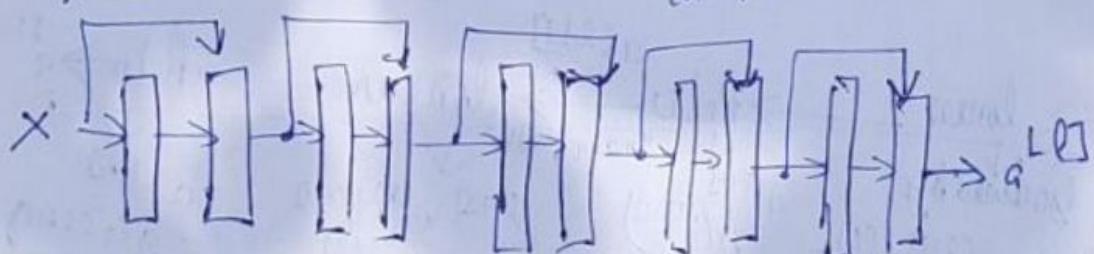
→ formula for this change to g = non-linearity

$$a^{[l+1]} = g(z^{[l+2]} + a^{[l]})$$

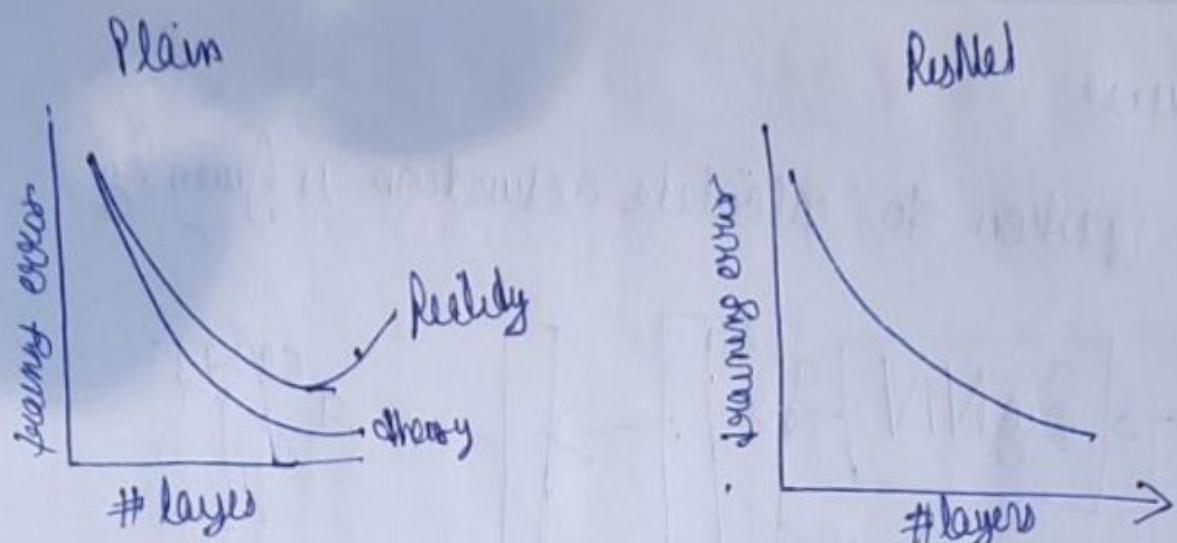
So, the addition of this $a^{[l]}$ here, it makes this a residual block

Residual Block Network

At has 5 residual block



And this is a residual network



- Performance of ResNet increases as the network goes deeper

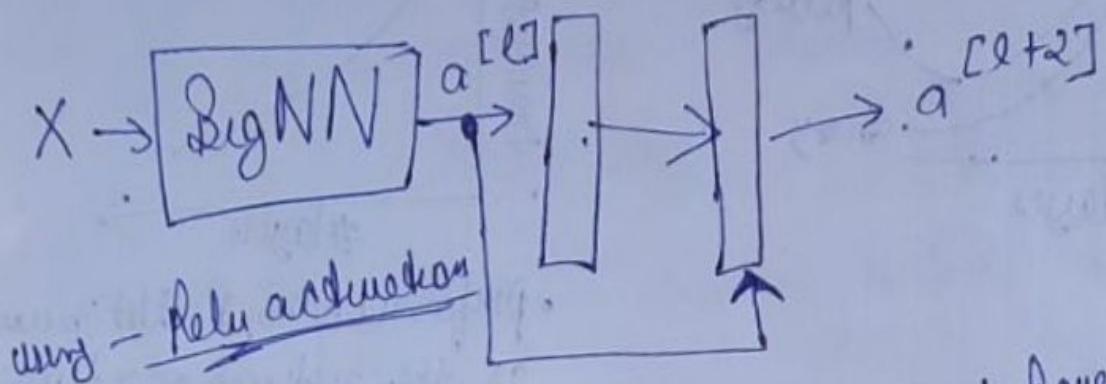
- In some cases going deeper won't effect the performance and that depends on the problem on your hand.
- Some people are trying to train 1,000 layer now which isn't used in practice

④ WHY RESNETS WORK

In order to make a good model, we first have to make sure that it's performance on the training data is good. We have seen earlier that training deeper networks using a plain network increases the training error after a point of time. But while training a residual network, this isn't the case. Even when we build a deep deep residual network, the training error generally does not

(increase)

The equation to calculate activation is given by



$$a^{[l+2]} = g(a^{[l+2]} + a^{[l]}) \rightarrow \begin{matrix} \text{they must have} \\ \text{same dimensions} \end{matrix}$$

$$= g(w^{[l+2]} a^{[l+1]} + b^{[l+2]} + a^{[l]})$$

Now, let's say $w^{[l+2]} = 0$ and the bias $b^{[l+2]}$ is also 0

then $a^{[l+2]} = g(a^{[l]})$

It is fairly easy to calculate $a^{[l+2]}$ knowing just the value of $a^{[l]}$.

970

⑤ NETWORKS IN NETWORKS AND 1X1 CONVOLUTIONS

Q

$$\begin{array}{|c|c|c|c|c|c|} \hline
 1 & 2 & 3 & 4 & 5 & 6 \\ \hline
 3 & 5 & 5 & 1 & 3 & 4 \\ \hline
 2 & 1 & 3 & 4 & 9 & 3 \\ \hline
 4 & 7 & 8 & 5 & 7 & 9 \\ \hline
 1 & 5 & 3 & 7 & 4 & 8 \\ \hline
 5 & 4 & 9 & 8 & 3 & 5 \\ \hline
 \end{array} \quad * \quad \boxed{2} \quad = \quad \begin{array}{|c|c|c|c|c|c|} \hline
 1 & & & & & \\ \hline
 2 & 4 & 6 & . & . & . \\ \hline
 & & & & & \\ \hline
 \end{array}$$

$6 \times 6 \times 1$

1x1 filter

But that's the case of $6 \times 6 \times 1$ channel image

If you have a $6 \times 6 \times 32$ image then a convolution with a 1×1 filter can do something that makes much more sense what a ~~one~~ 1×1 convolution will do is it will look at each of the 36 different position here (6×6) and it will take the element wise product b/w 32 numbers on the left and 32 numbers in the filter and then apply a ReLU non-linearity to it after that.

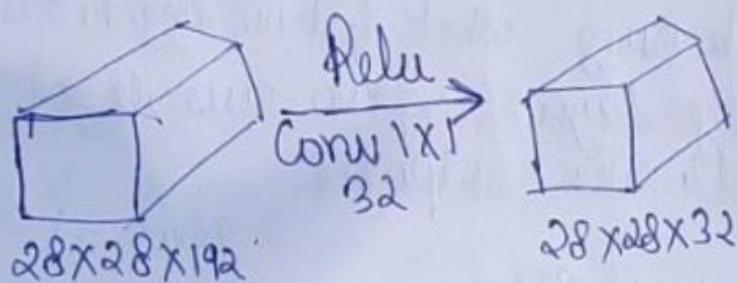
$$\begin{array}{ccc}
 \begin{array}{c} \text{one slice} \\ \text{---} \\ \text{6} \times 6 \times 32 \end{array} & * & \begin{array}{c} \text{---} \\ \text{1} \times 1 \times 32 \end{array} = \begin{array}{|c|c|c|c|c|c|} \hline
 1 & 1 & 1 & . & . & . \\ \hline
 & & & & & \\ \hline
 \end{array} \\
 & & \text{6} \times 6 \times \# \text{filters}
 \end{array}$$

So, to look on one of the 36 positions, maybe one slice through this value, you take these 36 numbers multiply it by / slice through the volume like that,

and we end up with a single number which then gets plotted in one of outputs like that.

- It is called 1×1 convolution but sometimes also called Networks in Networks.

Using 1×1 convolution

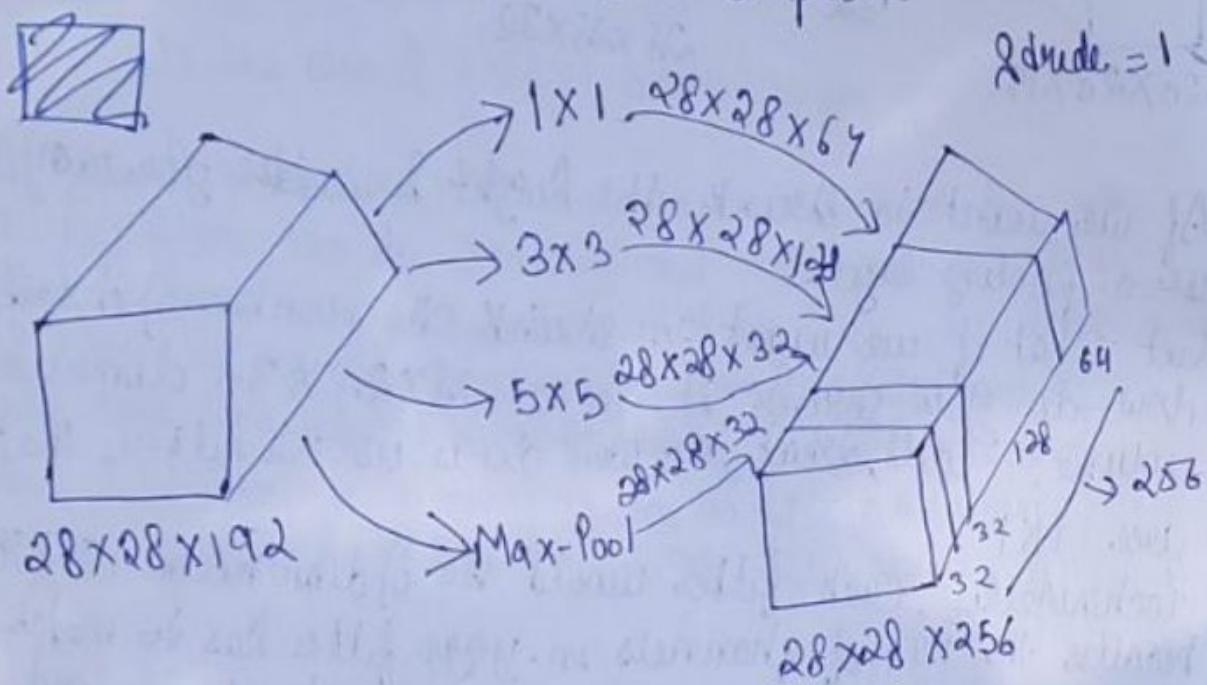


- If we want to shrink the height & width you can use a pooling layer
- But what if we want to shrink the number of channels how do you shrink it to a $28 \times 28 \times 32$ dimension volume? Well, what you can do is use 32 filters that are 1×1
- Technically, each filter would be 3 dimensions $1 \times 1 \times 192$ because the no. of channels in your filter has to match the number of channels in your input volume and the output of this process will be an $28 \times 28 \times 32$ volume
- This help in Building a inception network (next topic)

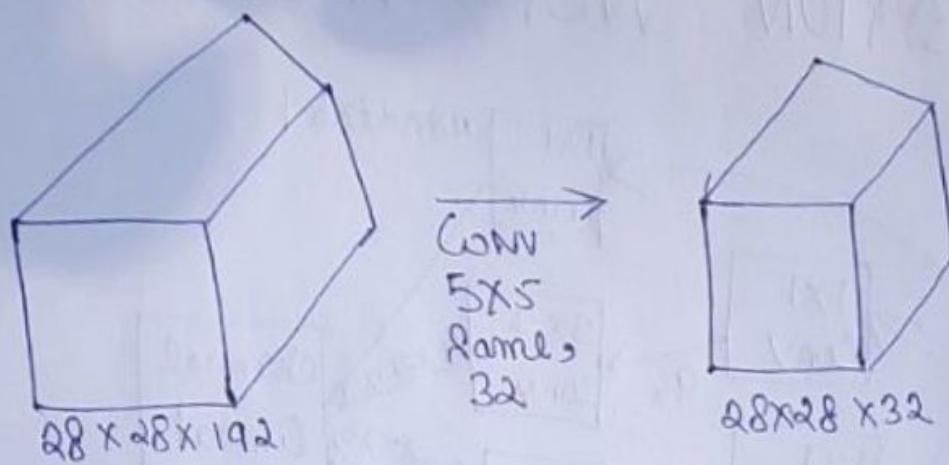
⑥ INCEPTION NETWORK MOTIVATION

- When you design a CNN you have to decide all the layers yourself. Will you pick a 3×3 conv or 5×5 conv or maybe a max pooling layer. You have so many choices.

Suppose we have $28 \times 28 \times 192$ input volume. Instead of trying to choose what filter size to use or conv layer or pooling layer. Inception uses all of them and stacks all the outputs.

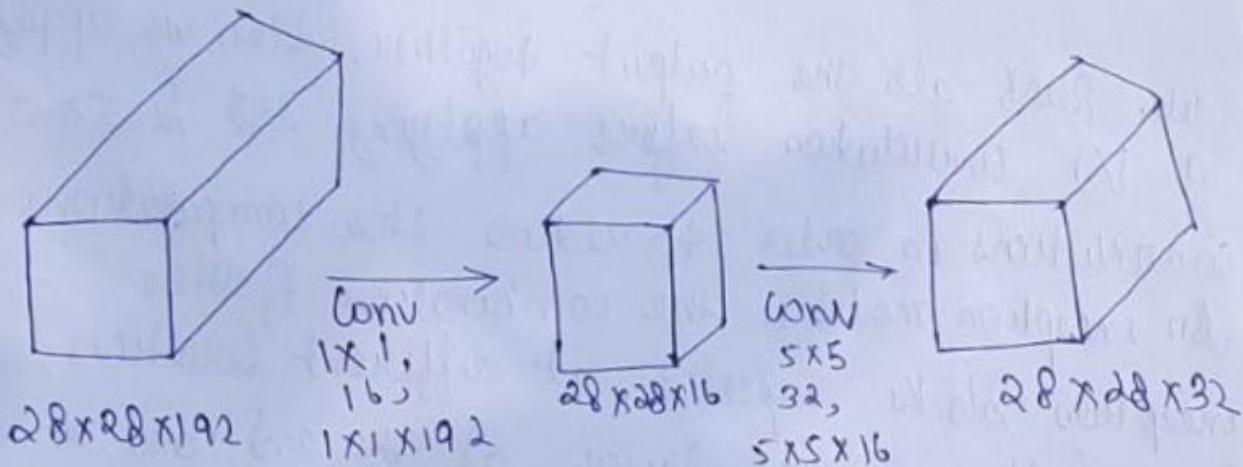


A good question to ask here - why are we using all these filters instead of using just a single filter size 5×5 ? Let's look at how many computations would arise if we would have used only a 5×5 filter on our input.



Number of parameters = $28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120\text{million}$
 & it has huge computation cost.

Now let's look at the computations at 1×1 convolution
 " and then a 5×5 convolution will give us:

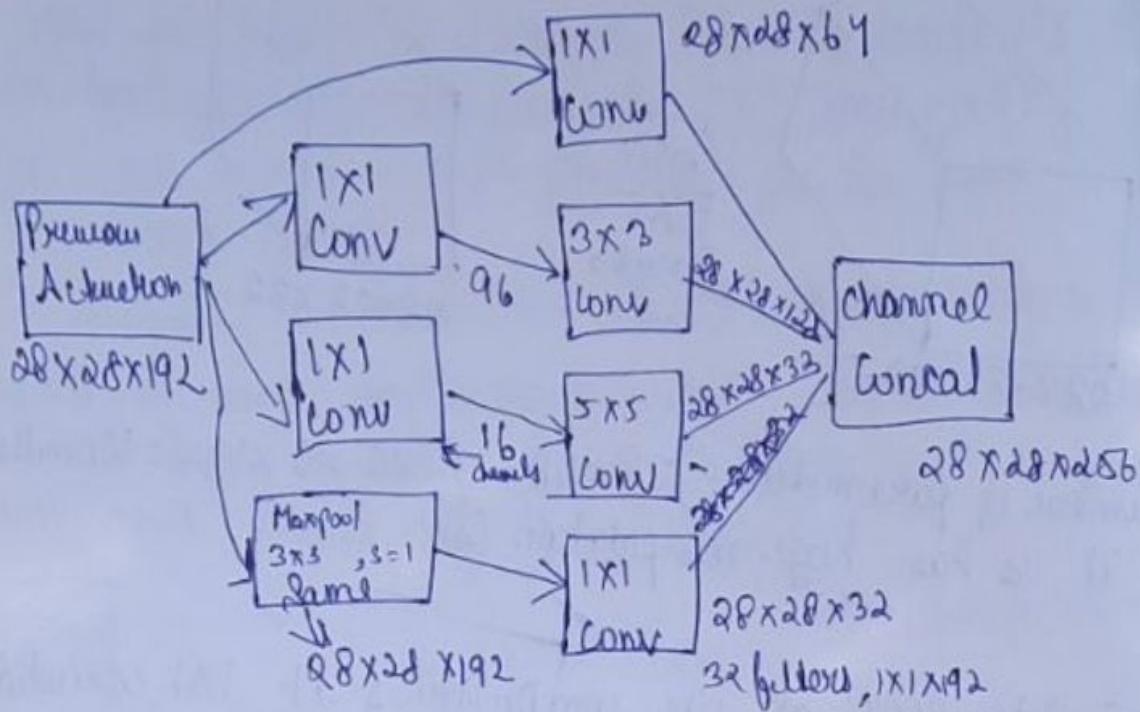


$$\begin{aligned} \text{No. of parameters in first convolution} &= 28 \times 28 \times 16 \times 1 \times 1 \times 192 \\ &= 2.4 \text{ million} \end{aligned}$$

$$\begin{aligned} \text{No. of parameters in 2nd " } &= 28 \times 28 \times 32 \times 5 \times 5 \times 16 \\ &= 10 \text{ million} \end{aligned}$$

Here we have seen a significant reduction. This is one key idea behind inception.

(⑦) INCEPTION NETWORK



We stack all the output together. Also, we apply a 1×1 convolution before applying 3×3 & 5×5 convolutions in order to reduce the computations. An inception model is the combination of these inception blocks repeated at different locations, some fully connected layer at the end, and a softmax classifier to output the classes.



PRACTICAL ADVICES FOR USING CONV-NETS

(1) Using Open-Source Implementation

- o It turns out that a lot of these NN are difficult to replicated because there are some details that may not presented on its papers. There are some other lessons like:
 - Learning decay
 - Parameter cleanup
- o A lot of deep learning researcher are opening sourcing their code into internet on sites like GitHub.
- o If you see a research paper and you want to build over it, the first thing you should do is to look for an open source implementation for this paper.
- o Some advantages of doing that you might download the network implementation along with its parameters/weights. The author might have used multiple GPUs and spent some weeks to reach this result and it might be front of you after you download it.

(2) Transfer learning

example:

- o Let's say you have a cat classification problem which contains 3 classes Tigger, Misty & neither

cat names

- You don't have much a lot of data to train a NN on these images
- Go online and download a good NN with its weights, the softmax function layer and put your own one & make the network learn only the new layer while other layer weight are frozen
- Frameworks have options to make the parameter frozen in some layers using trainable=0 or freeze=0
- One of the tricks that can speed up your training is to run the pretrained NN without final Softmax layer and get an intermediate representation of your images and save them to disk. And then use these representation to a shallow NN. This can save you the time needed to run an image through all the layers.
 - It's like converting your images into vectors
- Another example
 - what if in the last example you have a lot of pictures for your cats
 - one thing you can do is to freeze few layers from the beginning of the pretrained network and learn the other weights in the network.
 - Some other idea is to throw away the layers that aren't frozen and put your own layer there
- Another example

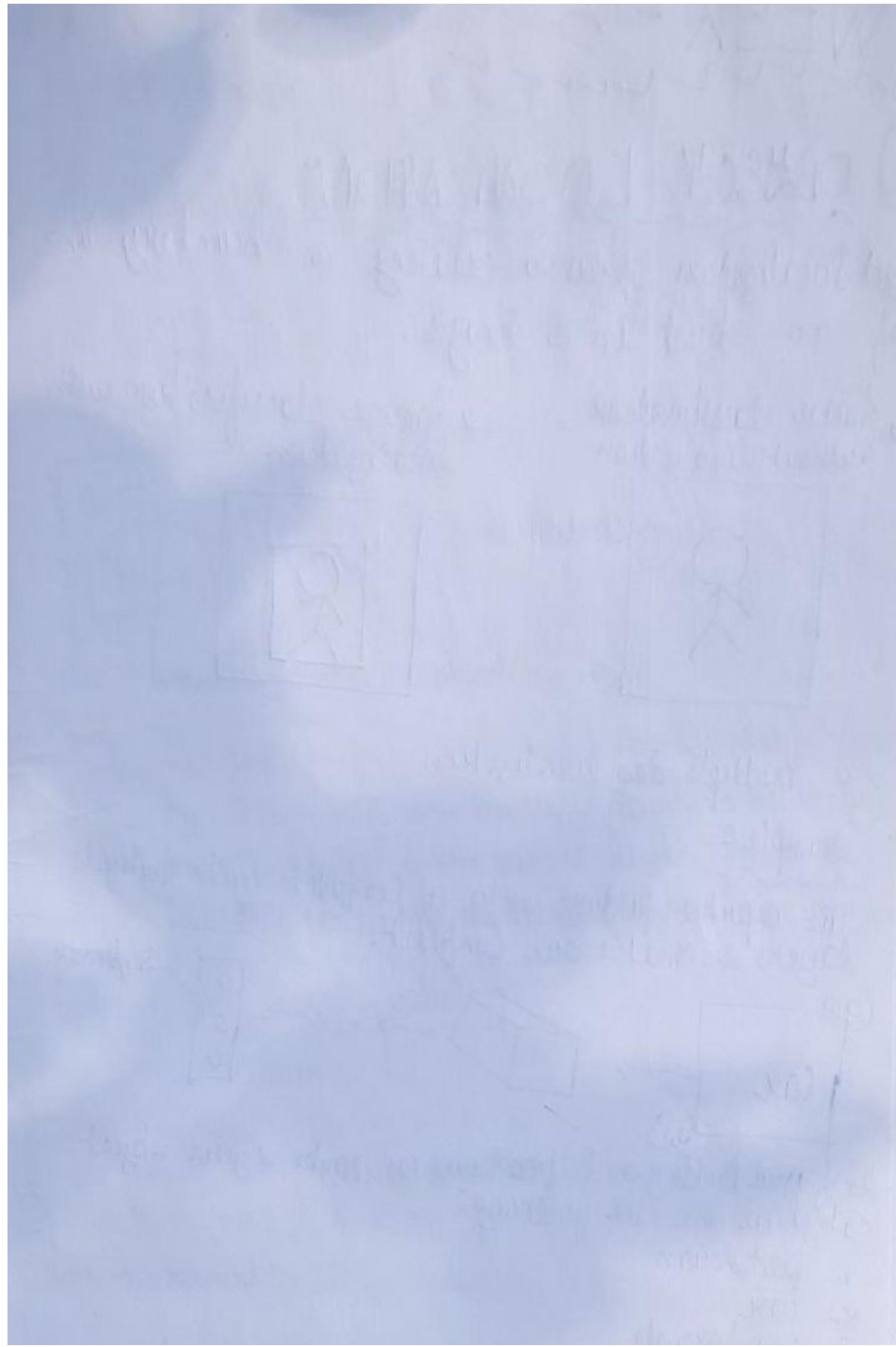
- If you have enough data, you can fine-tune all the layers in your pretrained network but don't randomly initialize the parameters, leave the learned parameters as it is and learn from them

③ DATA AUGMENTATION

Some

④ STATE OF COMPUTER VISION

from online Notes

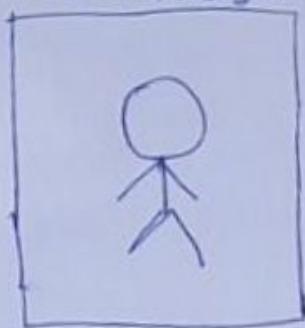


WEEK - 3

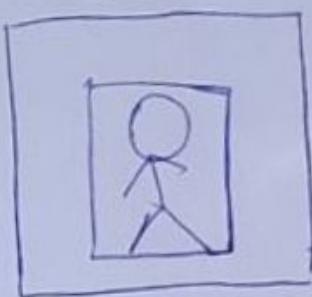
① OBJECT LOCALIZATION

object localization refers to creating a boundary line of an object in a image.

↳ Image classification without localization



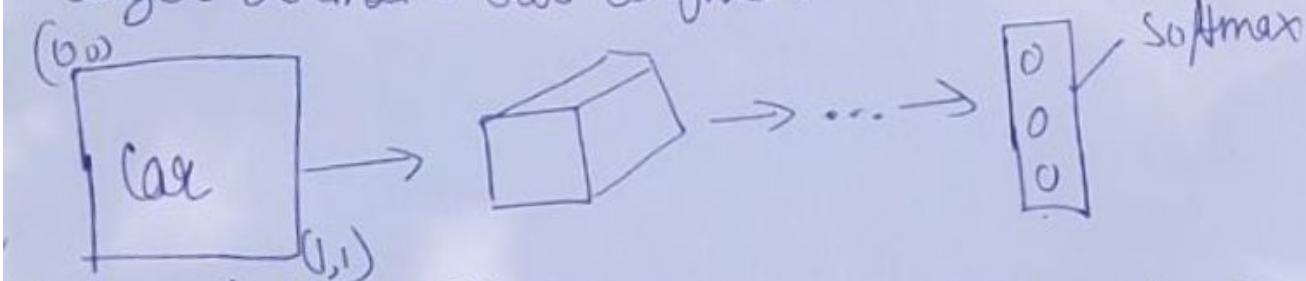
↳ Image classification with localization



* multiple localization

Example:

we input a picture into a ConfNet with multiple layers so that's our ConfNet.



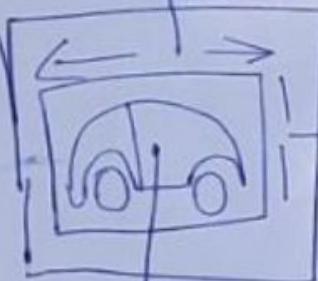
So if you're building a self-driving car maybe your object categories are the following:

1. pedestrian
2. car
3. motorcycle
4. background \rightarrow when there is nothing to detect

Now for for localizing the object we have to change
out N.N. we need to have few more output units

$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow 0$ Softmax

$\begin{bmatrix} bx \\ by \\ bw \\ bh \end{bmatrix}$ boundary box (bx, by, bw, bh)
 $\begin{array}{l} bx \\ by \\ bw \\ bh \end{array} = \begin{array}{l} \hookrightarrow \text{more out for localizing} \\ \hookrightarrow \text{additional output} \end{array}$


 bh (height of the boundary)

middle point of the object (bx, by)

$bx = 0.5$ (halfway to the right to the image)

$by = 0.7$ (70% to the way down to the image)

$bh = 0.3$ (30% of the overall height of the image)

$bw = 0.4$ (0.4 of the overall width of the entire image)

Defining the target label y

These are 4 classes

1. Pedestrian

2. Car

3. Motorcycle

4. background

} Need to output bx, by, bw, bh class label (1-4)

So, now let's define the target variable y

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_n \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

p_c - P_c = Probability that there is an object
 b_x, b_y, b_n, b_w → If there is an object then we want the location of that object.
 c_1, c_2, c_3 → It will tell us it class1, class2 or class3

If $p_c = 0$ then y will look like this

$$y = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

→ don't calculate

lets describe the loss function

$$L(\hat{y}, y) = \left\{ (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2 \right\} \text{ where are 8 output in } y$$

② LANDMARK DETECTION

example:

If you are working in a face recognition problem you might want ~~use~~ some points on the face like corners of the eyes, corners of the mouth, and corners of the

nose and so on. This can help in a lot of application like detecting the pose of the face.

- Suppose we need 64 landmarks then Y will have 65 output

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{64} \end{bmatrix} \rightarrow \text{As there is face if yes then } Y = \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \rightarrow \text{No face}$$

③ FACE DETECTION

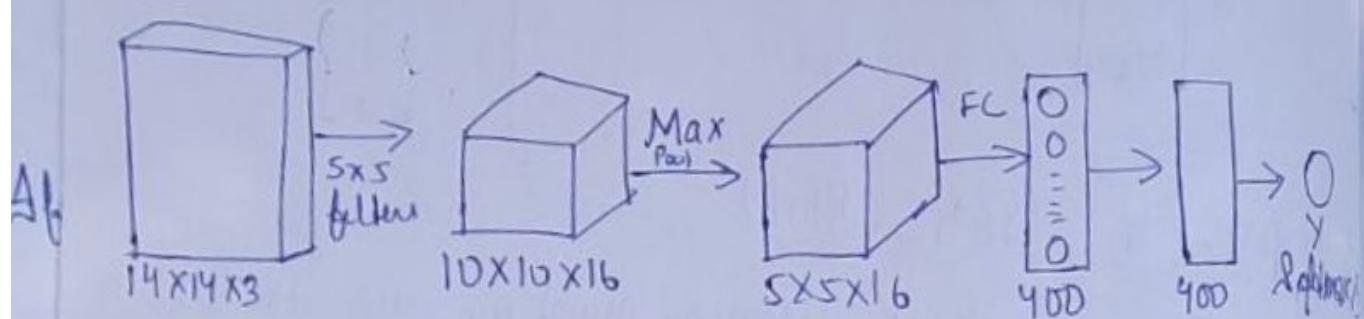
- To solve the object detection problem we use a technique called sliding window detection.
- We move the sliding window on the image to detect the image and moving that sliding window box ~~is~~ on ~~different~~ image then we zoom the image little bit then we move it again other again same process.
- Disadvantage of sliding window is computation time.
- To solve this problem, we can implement the sliding window with a convolutional approach.

④ CONVOLUTIONAL IMPLEMENTATION OF SLIDING WINDOWS

So, 'to Another' we'll learn how do implement object detection of convolutionally do make this algo of work fast

Turning FC layer into convolutional layers

So, let's say that your object detection algorithm input
 $14 \times 14 \times 3$ images



Now we'll see how these layers turned into convolutional layers

from video and notes

⑤ BOUNDING BOX PREDICTIONS

In this we'll learn how to make your predictions to be more accurate

Sometimes none of the sliding window boxes really match up perfectly with the position of the car

A good way to get the output more accurate bounding boxes is with the YOLO algorithm.

- Yolo stands for You only look once
- Let's say we have a image I 100×100
- We use 3×3 grid on the images this will divide our image in a 9 parts like this

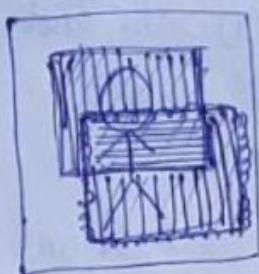


But

In Yolo algorithm we'll have problem if we have found more than one object in one grid. Yolo ~~is~~ algorithm has a great speed and a low net implementation. Yolo uses a single CNN network for both classification & localizing the object using bounding boxes.

⑥ INTERSECTION OVER UNION

It computes the intersection over union of these two bounding boxes



Intersection over Union (IoU)

$$= \frac{\text{Size of } \cap}{\text{Size of } \cup}$$

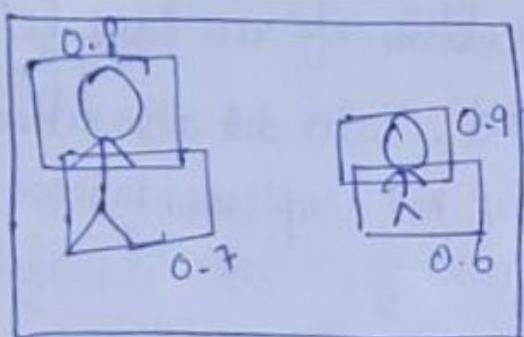
"Correct" if $\text{IoU} \geq 0.5$

We can change the 0.5 , 0.6 or 0.7 or any value.

IoU measure of the overlap b/w two bounding boxes
the higher the IoU the better is the accuracy

⑦ NON-MAX SUPPRESSION

- One of the problem we face with Yolo algorithm
 - it can detect an object multiple times
- Non-max suppression is used so that Yolo algorithm detect object just one time.



It will share the rectangle with highest Probability

Non-max sup. algo.

→ Y shape should be $[P_c, bx, by, bn, bw]$

→ discard all boxes with $P_c < 0.6$

→ Pick the box with largest P_c

→ discard any remaining box with $\text{IoU} > 0.5$ with \downarrow

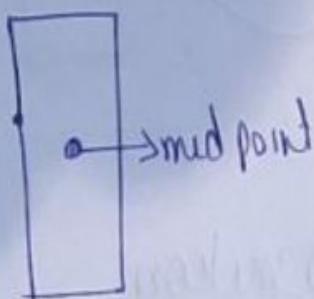
⑧ ANCHOR BOXES

watch video

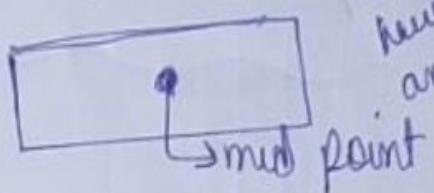
In yolo algorithm, a grid only detects one object.
What if a grid cell wants to detect multiple objects?

In this situation we use Anchor Box

Anchor Box 1



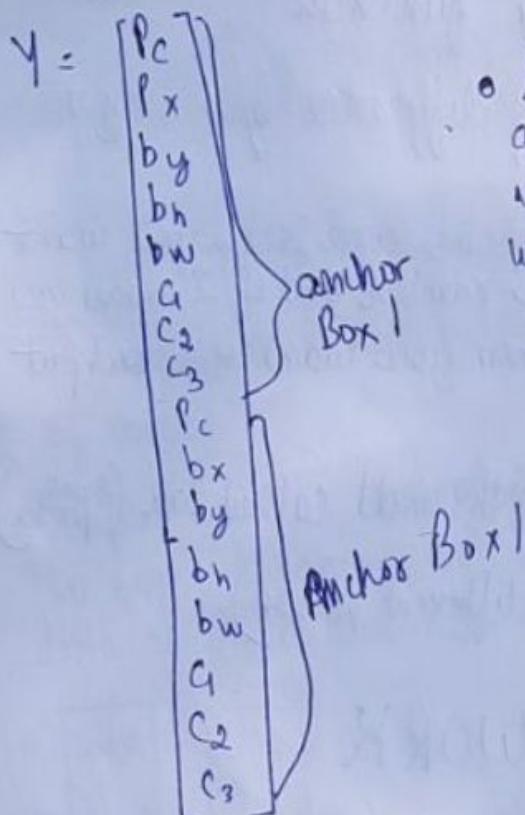
Anchor Box 2



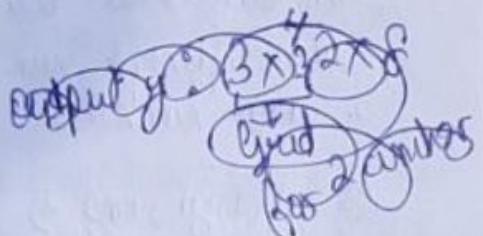
We generally have 4 or 5 anchor boxes.

With 2 anchor box, Each object in drawing image is assigned to grid cell that contains object's midpoint & anchor box for the grid cell with highest IOU. You have to check where your object should be based on its rectangle closest to which anchor box.

Example :-



- Anchor Boxes allows your algorithm to specialize, means in our case to easily detect wider images or taller ones.



With wider -

⑨ YOLO ALGORITHM

WEEK - 4

① What is face recognition

Nothing to learn

② ONE SHOT LEARNING

- It refers to when a recognition system is able to recognize a person face fo just from one image of a person

Learning a "similarity" function

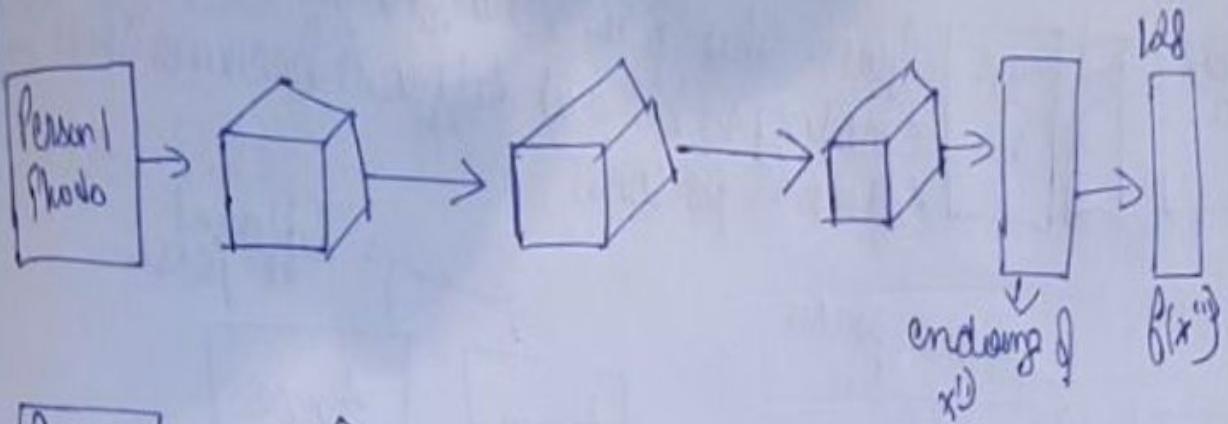
$$d(\text{img1}, \text{img2}) = \text{degree of difference b/w images}$$

So if the two images are of the same person, you want this output to be a small number and if 2 images are two different people then you want it to output a large number.

If $d(\text{img1}, \text{img2}) \leq \tau \rightarrow$ threshold called τ (same person)

$d(\text{img1}, \text{img2}) > \tau$ (different person)

③ SIAMESE NETWORK



$$d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_2^2$$

So, this idea of running two identical, convolutional NN works on two different inputs and then comparing them, it is called Siamese NN architecture.
 These 2 NN ~~has~~ has same parameters.

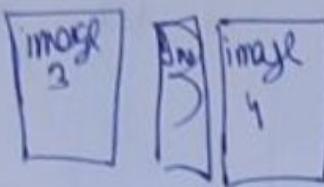
④ TRIPLEX LOSS

To apply the triplet loss, you need to compare pairs of images

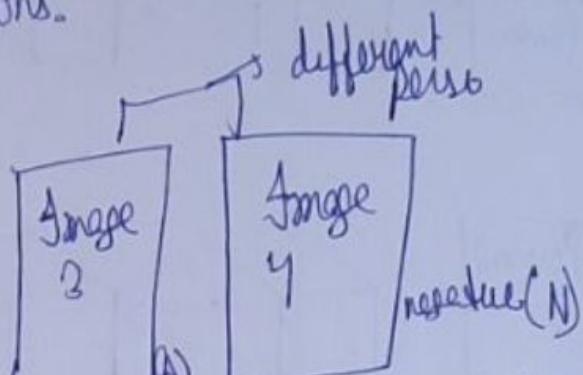
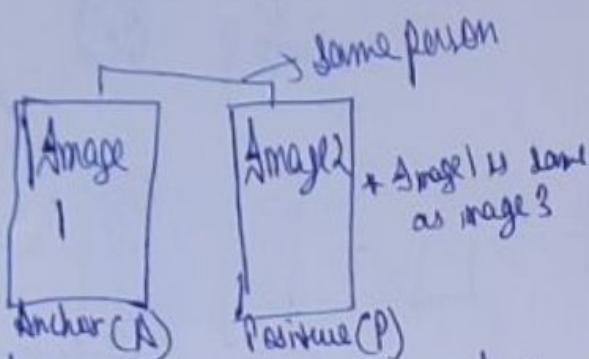
Triplet loss is one of the loss functions we can use to loss functions we can use to solve the similarity distance in a Siamese network.



Given this pair of images, you want their encoding to be similar because these are the same person.



Given this pair of images you want their encoding to be quite different because these are different persons.



(As the terminology of the triplet loss, always look at one anchor image and then you want to distance b/w the anchor and the positive image)

* The difference b/w img1 & 2 should be less and the difference b/w img3 and 4 should be more this is what gives rise to the term triplet loss which means you'll always looking at three images at a time (anchor, positive and negative img at time)

So to formulate this what you want is for the parameters of your N.N. of your encoding to have following property

$$\underbrace{\|f(A) - f(P)\|^2}_{\text{which is } d(AP)} \leq \underbrace{\|f(A) - f(N)\|^2}_{\text{which is } d(AN)}$$

Now if we move the equation

$$\|f(A) - f(P)\|^2 - \|f(A) - f(P)\|^2 \leq 0$$

But we are going to change the equation a little which is one small way to make sure this is satisfied, if everything equals zero

So, to make sure that the NN doesn't just output zero, for all the encodes. So to prevent this what we are going to do is modify (It needs to be quite a bit smaller than zero)

so, we add margin (α) here.

$$\|F(A) - F(P)\|^2 - \|F(A) - F(N)\|^2 + \alpha \leq 0$$

\downarrow
alpha or
margin

Triplet loss function

Given 3 picture A, P, N

$$\ell(A, P, N) = \max\left(\|F(A) - F(P)\|^2 - \|F(A) - F(N)\|^2 + \alpha, 0\right)$$

Triplet Cost function

$$y = \sum_{i=1}^m \ell(A^{(i)}, P^{(i)}, N^{(i)})$$

~~Suppose you have to training 10,000 pictures with 1,000 different persons~~

Notice that in order to define this dataset of triplets you do need some pairs of $A \neq P$. Pairs of pictures of the same person.

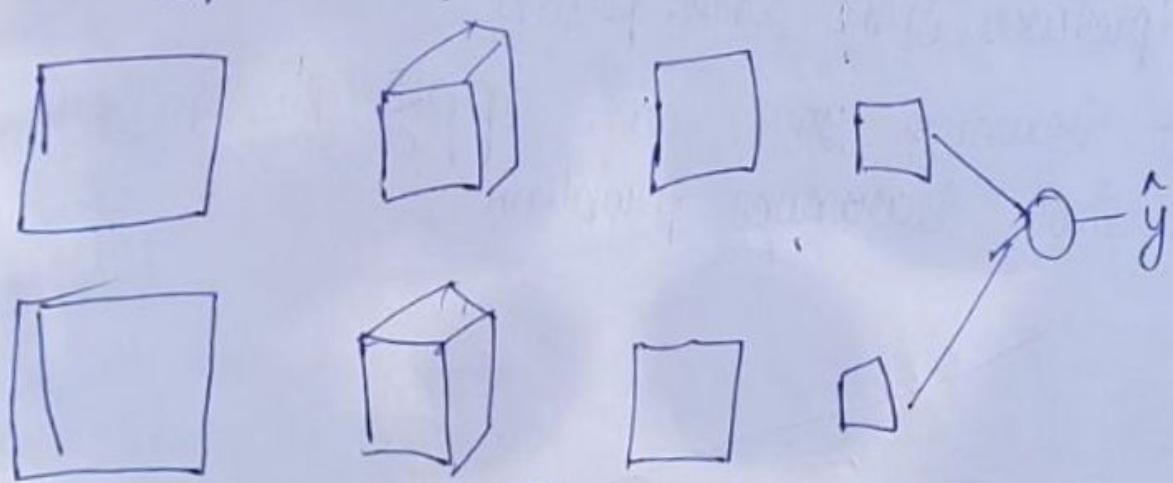
After training you can apply this to your one shot learning problem

⑥ FACE VERIFICATION AND BINARY

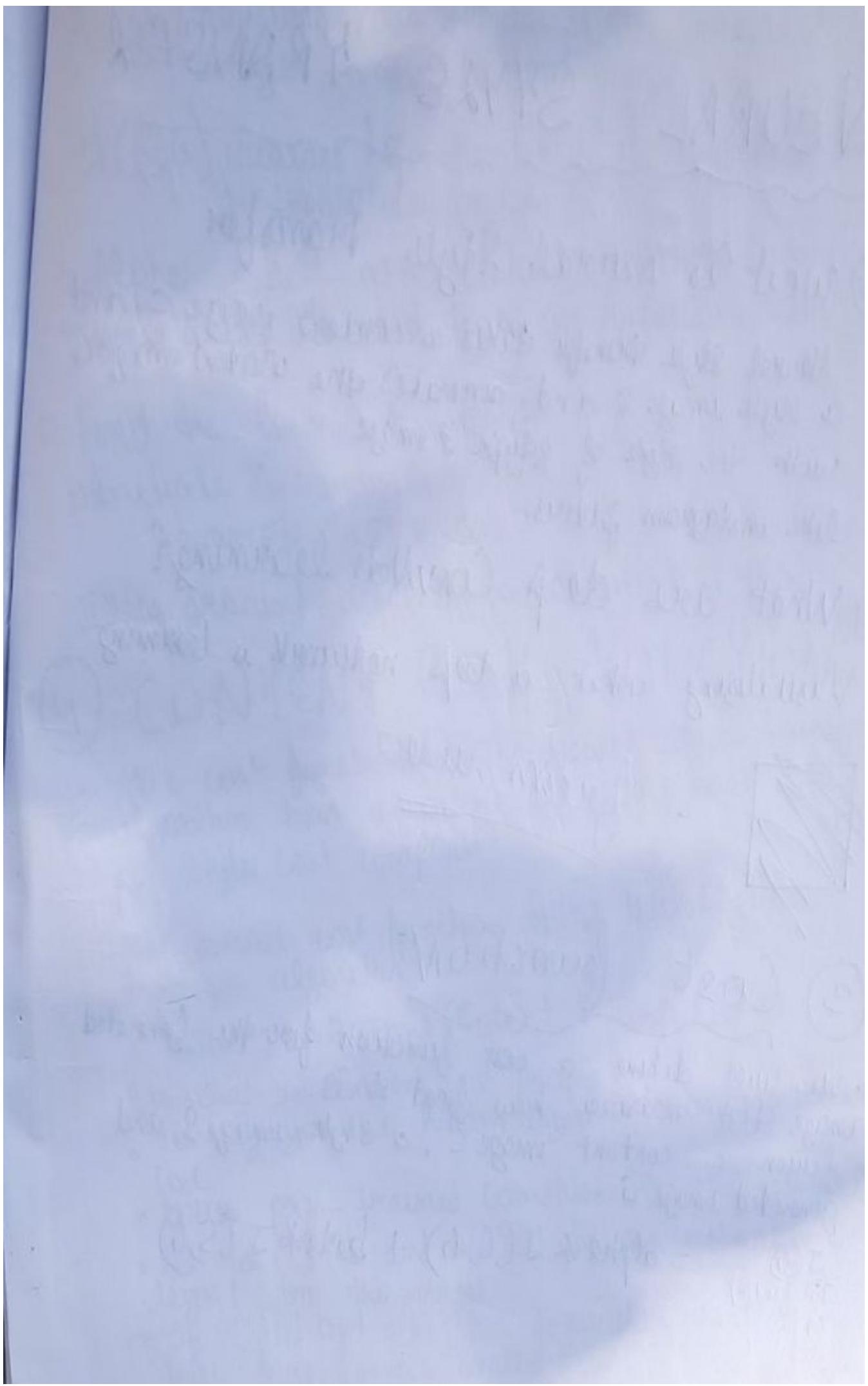
~~Face recognition can also be posed as a straight binary classification problem~~

The triplet loss is one good way to learn the parameters of a confnet for face recognition. There another way to learn these parameter let me show you how face recognition can also be posed as a straight binary classification problem

Another way to obtain a NN is to take pair of NN (Siamese Network) and have them compute the embeddings



From notes online



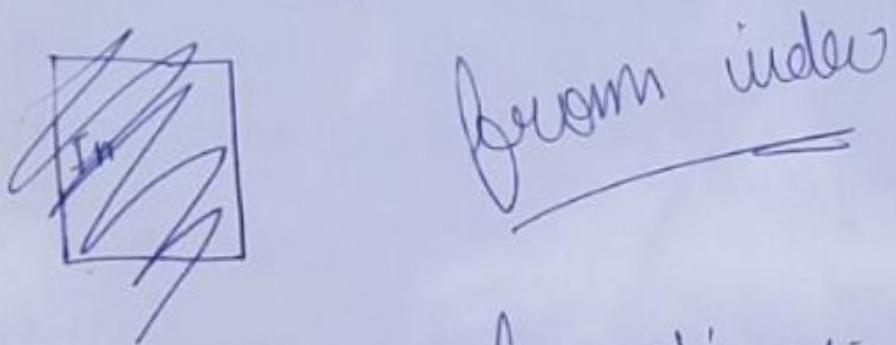
NEURAL STYLE TRANSFER

① What is neural style transfer

Neural style transfer takes a content image C and a style image S and generates the content image G with the style of style S image like Instagram filters.

② What are deep ConvNets learning?

Visualizing what a deep network is learning



③ Cost function

- we will define a cost function for the generated image that measures how good it is.
- Given a content image C , a style image S and generated image G .

$$J(G) = \text{alpha} * J(C, G) + \text{Beta} * J(S, G)$$

$J(G)$
cost func of
 C

$\mathcal{L}(C, G)$ measures how similar is the generated image to the content image

$\mathcal{L}(S, G)$ measures how similar is the generated image to the style image

- Alpha & Beta are relative weighting do the similarity \rightarrow and these are hyperparameters.

i) find the generated image G

ii) mutate G randomly

o example : $100 \times 100 \times 3$

iii) use gradient descent to minimize $\mathcal{L}(G)$

④ CONTENT COST FUNCTION

The cost function of the Neural style transfer algorithm had a content cost component and a style cost component.

The overall cost function of the Neural style transfer algorithm

$$J(G) = \alpha \mathcal{L}_{\text{content}}(G, G) + \beta \mathcal{L}_{\text{style}}(S, G)$$

So what is content cost function?

• lets say you use hidden layer 1 to compute content cost.

• use pre-trained ConvNet (Eg VGG network)

• let $a(C)[1]$ and $a(G)[1]$ be the activation of layer 1 on the image

• If $a(C)[1]$ and $a(G)[1]$ are similar then they will have same content

$a(z, t)$ at a layer $l = \frac{1}{2}$

$$J_{\text{vertical}}(z, t) = \frac{1}{2} \| a^{[l+1](t)} - a^{[l](t)} - a^{[l-1](t)} \|^2$$

⑤ SHYKE COSH FUNCTION



#5 SEQUENCE MODELS

WEEK - 1

① WHY SEQUENCE MODELS

examples of sequence models are :

- ↳ speech recognition (from audio to text)
- ↳ Music generation
- ↳ sentiment classification
- ↳ DNA sequence analysis
- ↳ Machine translation (one language to other)
- ↳ Video activity recognition
- ↳ Name entity recognition (find person name from sentence)

② NOTATION

x : Harry Potter and Hermione Granger invented a new spell.

lets say we want to find the names in this sentence (name entity recognition)

Now we need a model to output y which has one output per input word

y : 1 1 0 1 1 0 0 0 0 (sequence of 9 words)

Now, the input is the sequence of nine words. So, eventually we're going to have nine set of features to represent these 9 words.

X : Harry Potter and Hermione Granger invented a new spell
 $x^{(1)} \quad x^{(2)} \quad x^{(3)} \quad \dots \quad x^{(t)} \quad \dots \quad x^{(9)}$

We use $x^{(t)}$ to index into the positions in the sequence.
 $T_x = 9$ (length of input)

y :
 $y^{(1)} \quad y^{(2)} \quad y^{(3)} \quad 0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 0 \quad 0$
 $y^{(4)} \quad y^{(5)} \quad y^{(6)} \quad y^{(7)} \quad y^{(8)} \quad y^{(9)}$
 $T_y = 9$.

Note: T_x & T_y can be different

$x_i \quad x^{(i)}$ = i-th training example

$x^{(1)} \dots x^{(T_x)}$ = is the sequence of training example.

$T_x^{(i)}$ = input sequence length for training example i

Similarly $y^{(i)}, y^{(i)}, T_y^{(i)}$.

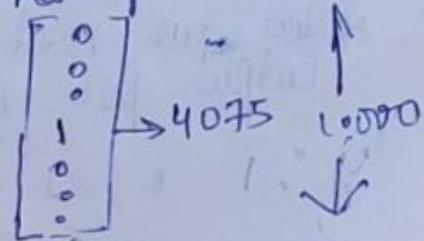
To represent a word in the sentence the first thing you do is come up with a vocabulary (called dictionary)

Vocabulary

a	1
aron	2
and	3 6 7
-	-
Harry	4 0 7 5
Potter	6 8 3 0
-	-
zuler	1 0 0 0 0 or any no.

what you can do is then use one hot representations to represent each of these words. For example the representation for Harry will be like

0s →
Same for all four words



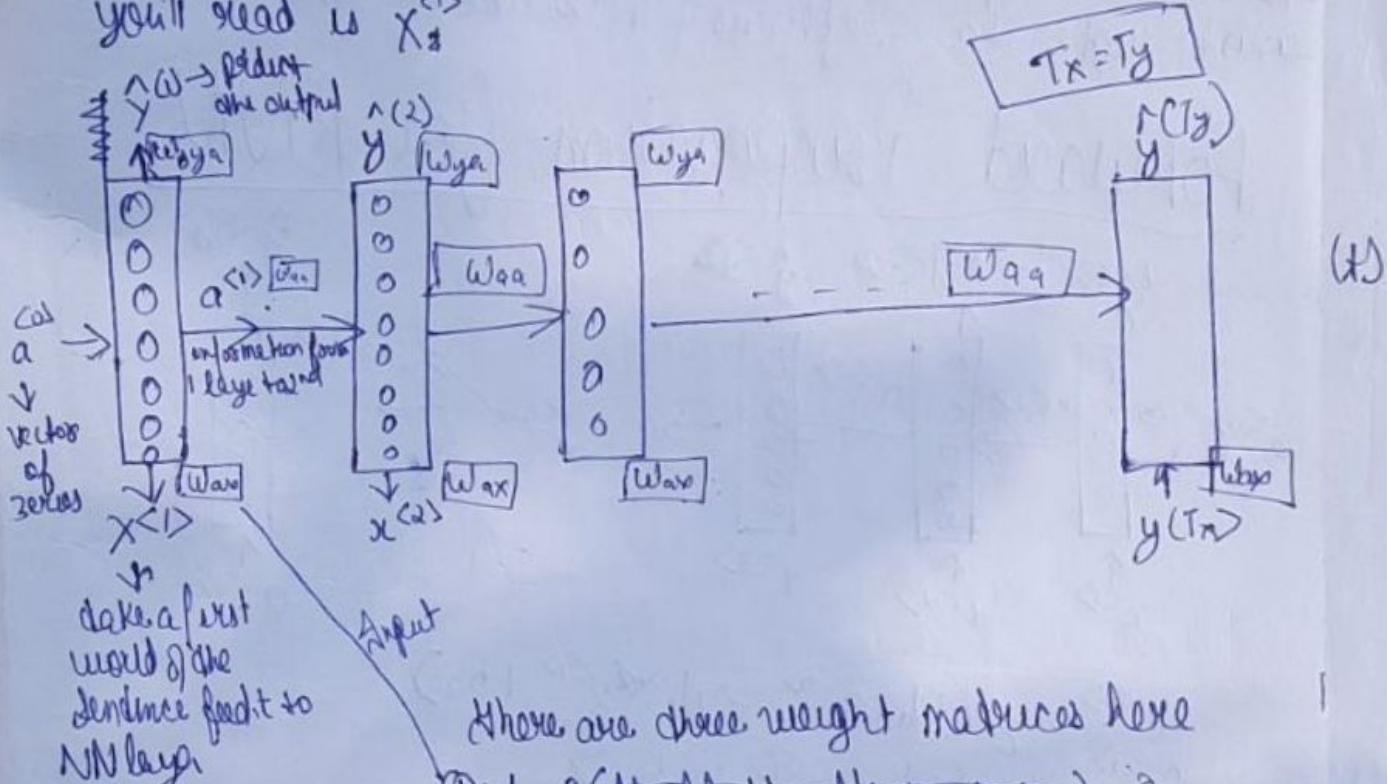
③ RECURRENT NEURAL NETWORK

The thing you could do is try to use a standard NN for this task. But it will have problem then:

- Input, outputs can be different lengths in different examples
- Doesn't share features learned across different position of text.

So, RNN don't have these advantages.

∴ if you are reading from left to right then the first you'll need to $x^{(1)}$



There are three weight matrices here

① W_{ax} : (~~No. of hidden neurons~~) \times $(No. of input neurons)$

② W_{ya} : (~~No. of hidden neurons~~, ~~No. of hidden neurons~~) \times $(No. of output neurons)$

③ W_{cy} : (~~No. of hidden neurons~~) \times $(No. of output neurons)$

↳ output predictions.

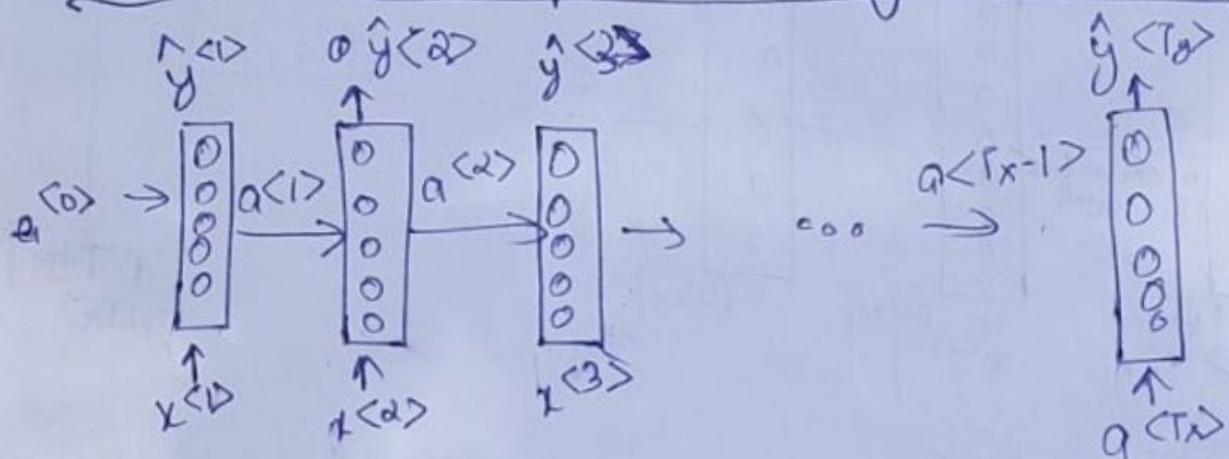
For prediction of $\hat{y}^{<3>}$ it gets information from x^1, x^2 and x^3

One weakness of RNN is that it only uses the information that is earlier in the sequence to make a prediction (left to right). for our sentence. If we can detect Harry is a name then it will output at $t=1$, but full name is Harry Potter so that's the problem.

$x: \underline{\text{Harry}} \text{ Potter & Hermione Granger invented a new spell.}$

~~An example~~ So we also need to go from right to left we'll talk about this later

Forward Propagation for RNN



$$a^{<1>} = g(W_a a^{<0>} + W_x x^{<0>} + b_a)$$

$$\hat{y}^{<1>} = g(W_y a^{<1>} + b_y)$$

~~but~~ activation function are often Tanh or ReLU

$$a^{<t>} = g(W_a a^{<t-1>} + W_x x^{<t>} + b_a)$$

$$\hat{y}^{(t)} = g(W_y a^{<t>} + b_y)$$

$$a^{(t)} = g(w_a [a^{(t-1)}, x^{(t)}] + b_a)$$

$$w_a = \begin{bmatrix} w_{aa} & | & w_{ax} \\ \xleftarrow{100} & & \xrightarrow{10,000} \\ \xleftarrow{100} & \text{stacking together} & \end{bmatrix}$$

$$\boxed{\begin{array}{l} \text{if } w_{aa} = (100, 000) \\ w_{ax} = (100, 10, 000) \end{array}}$$

$$\text{then } w_a = (100, 10, 100)$$

and what does mean $[a^{(t-1)}, x^{(t)}]$?

just take 2 vectors & stack them together

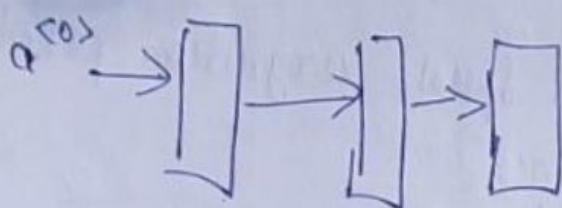
$$[a^{(t-1)}, x^{(t)}] = \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} \quad \begin{array}{c} \uparrow 100 \\ \downarrow 100 \\ \uparrow 10000 \\ \downarrow 10000 \end{array} \quad \begin{array}{c} \uparrow \\ \uparrow 10100 \\ \downarrow 10100 \end{array}$$

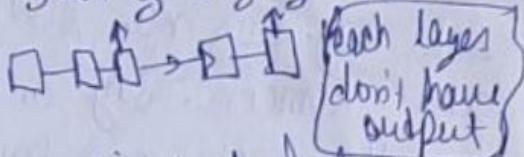
$$[w_{aa} \mid w_{ax}] \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} = w_{aa} a^{(t-1)} + w_{ax} x^{(t)}$$

④ BACKPROPAGATION THROUGH TIME

⑤ DIFFERENT TYPES OF RNNs

① Many to Many



* In this no. of input = no. of output
But sometimes they are not equal eg language translation


This architecture we have discussed before

② Many to one

Used in Sentimental Analysis



The output will be 1-5 like stars in Reviews

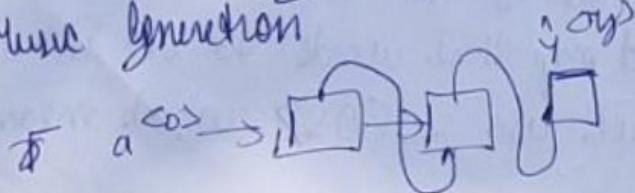
③ One to one

This is like standard N.N



④ One to many many

Music Generation



⑥ LANGUAGE MODEL AND SEQUENCE GENERATION

What is language model

Let's say we are solving a speech recognition problem and someone say a sentence

- The apple and pear deled
- The apple and pear deled.

Here pear & pair sounds exactly the same so how would a speech recognition application choose from the two.

That's where the language model comes in let see how.

$$P(\text{The apple I pair deled}) = 3.2 \times 10^{-13} \quad \text{with these probabilities the 2nd sentence is much more likely by over a factor of 10 to 100 compared to the 1st sentence.}$$
$$P(\text{The apple I pear deled}) = 5.7 \times 10^{-10}$$

What the language model does is it estimates the probability of that particular sequence of words.

$$P(y^{(1)}, y^{(2)}, \dots, y^{(T_y)})$$

To train RNN we need training set.

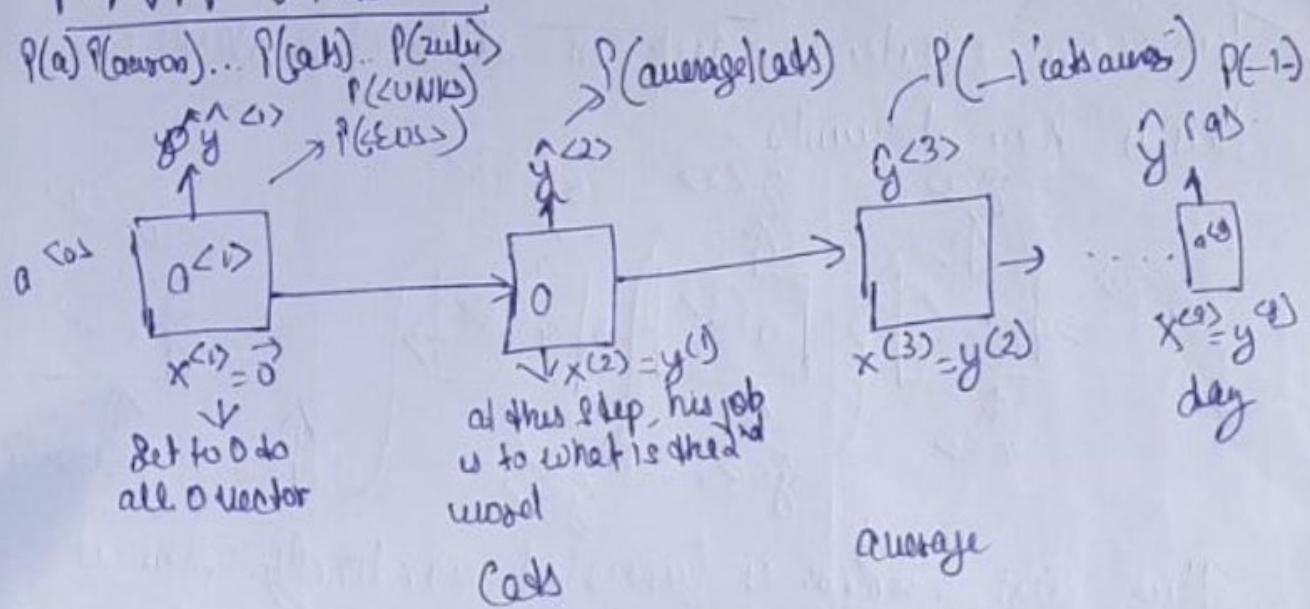
Training set : large corpus of english text

Suppose we have a sentence \rightarrow Gets average 15 hour of sleep a day.

The first thing will do is tokenize this sentence which means we'll form a vocabulary and map those words to one hot vector at the end of the sentence we use $\langle \text{EOS} \rangle$ which means end of sentence

So what if our sentence has a word which is not in our vocabulary then we assign those word a unique token UNK (Unknown words)

RNN MODEL



→ cars average 15 hours & sleep a day. <EOS>

and so on this RNN learn to predict one word at a time going from left to right.

The loss function is defined by cross-entropy loss.

$$L(y^{(t)}, \hat{y}^{(t)}) = - \sum_i y_i^{(t)} \log \hat{y}_i^{(t)}$$

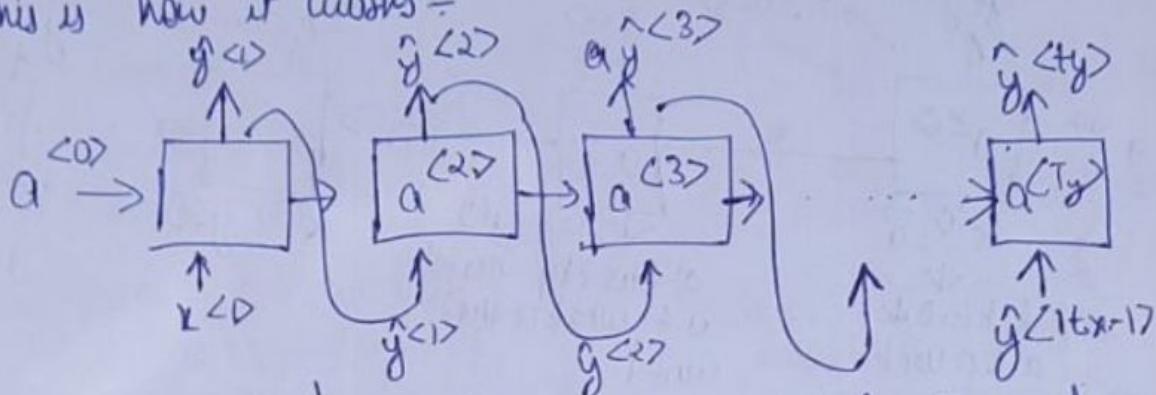
i is for all elements in the corpus,

t - for all timesteps

⑦ SAMPLING NOVEL SEQUENCES

It refers to getting to sample words from a trained model of RNN. like you trained a model of Shakespeare's text and you want text out of this model. You can create a text in style of Shakespeare.

This is how it works:



The first word is generated randomly then it finds the probability of next word the probable word with maximum probability is selected (from the vocabulary). This is how it is done.

It will keep going until we get the $\langle \text{EOS} \rangle$ token. we can reject any $\langle \text{UNK} \rangle$ token.

⑧ VANISHING GRADIENTS WITH RNN.

(6)

2018

(6) МИЛЫАДАСА АСЫРДО

① GATED RECURRENT UNIT (GRU)



⑩ LONG SHORT-TERM MEMORY (LSTM)

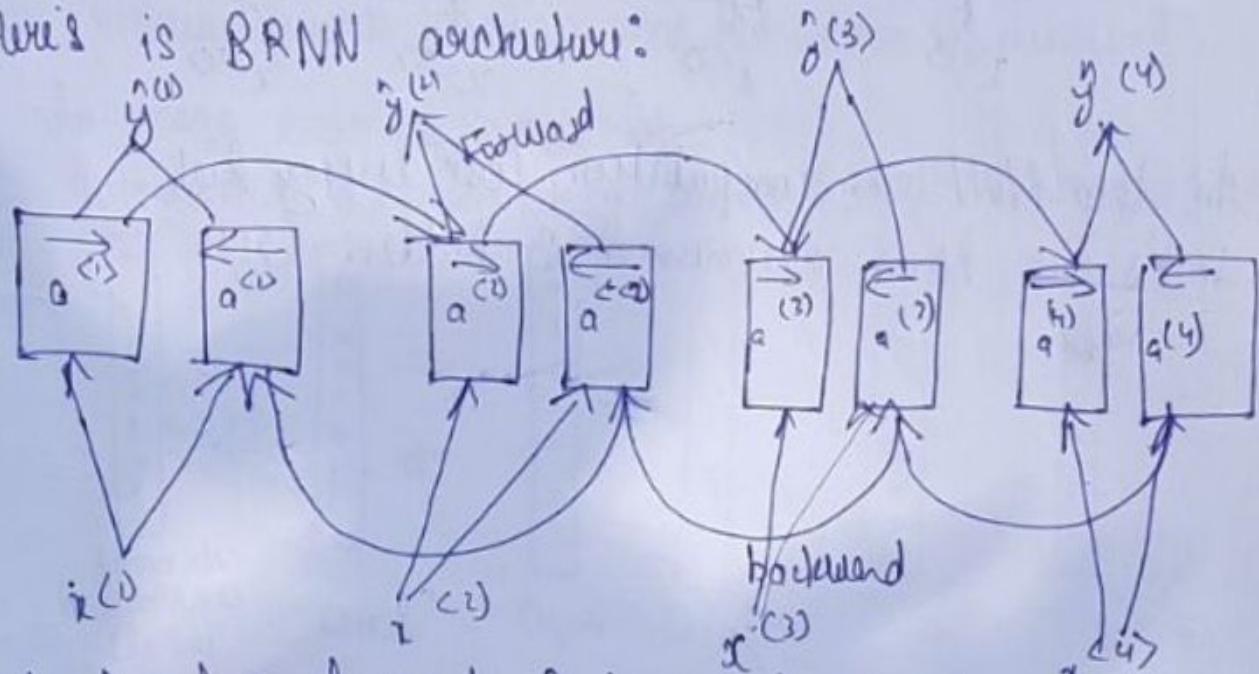
ДОКУМЕНТЫ ПОДДЕРЖАНИЯ

11) BIDIRECTIONAL RNN

What does it do?

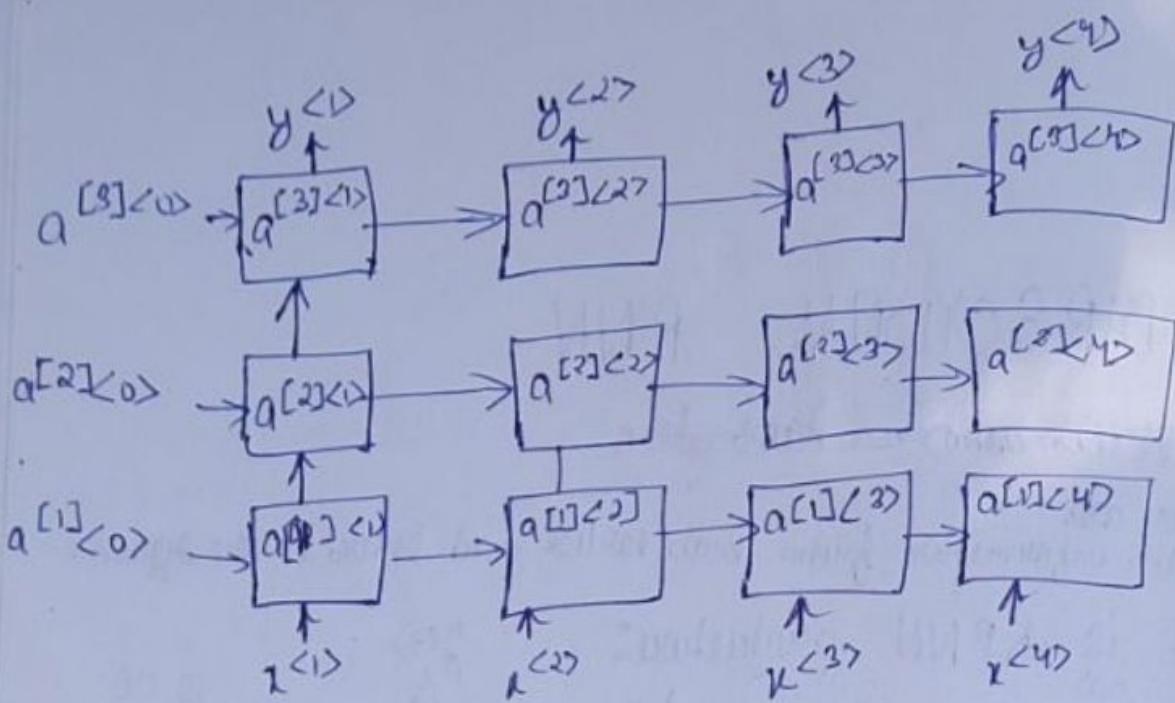
At take information from both earlier and later in the sequence

Here is BRNN architecture:



- At has both forward & backward info. Going x
- Note, that BiRNN is an cyclic graph.
- The block here can be RNN, LSTM or GRU
- In live Speech recognition if you use BiRNNs you will need to wait for the person who speaks to stop to take the entire sequence. and then make your predictions

⑫ DEEP RNNs



In deep RNN the computation cost is very high
that's we have maximum of 3 hidden layer

~~cost~~

WEEK - 2

[NLP is revolutionized by DL]

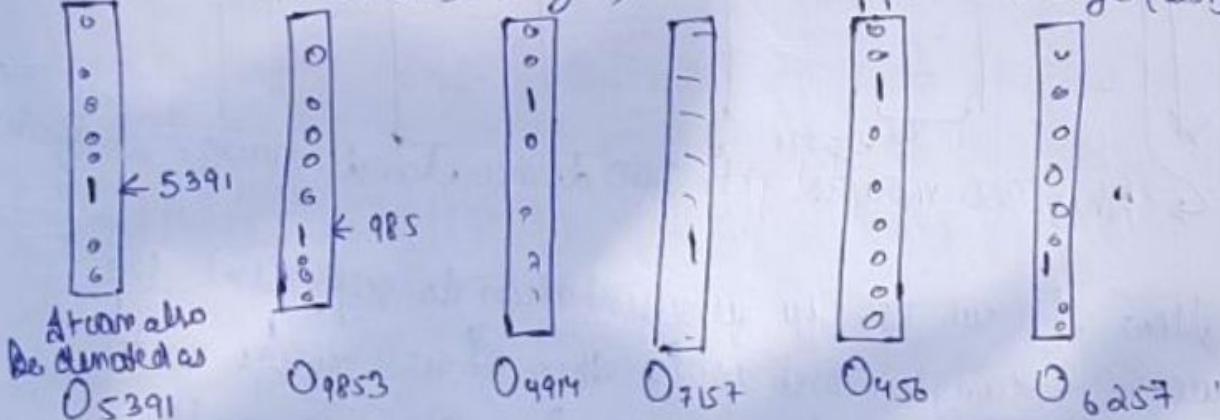
① WORD REPRESENTATION

Word embeddings which is a way of representing words
An other way algorithms automatically understand analogies
like : man is to woman, as King is to Queen

We represent words using a one hot vector as discussed.

An image example would be.

Man (sgn) Woman (q853) King (q914) Queen (q157) Apple (q56) Orange (q57)



One of the weakness of this representation is that it treat a word as a string that itself and it doesn't allow an algorithm to generalize across words.

Example : "I want a glass of orange ____". (a model will predict that next word is juice)

But what if we change "orange" to "apple"? Then a model can't predict the next word "juice" because it was trained orange not apple.

So, even it doesn't know that somehow apple and orange are much similar than King & orange or green or orange

So, instead of one-hot presentation, we can use featured representation.

WORD EMBEDDING

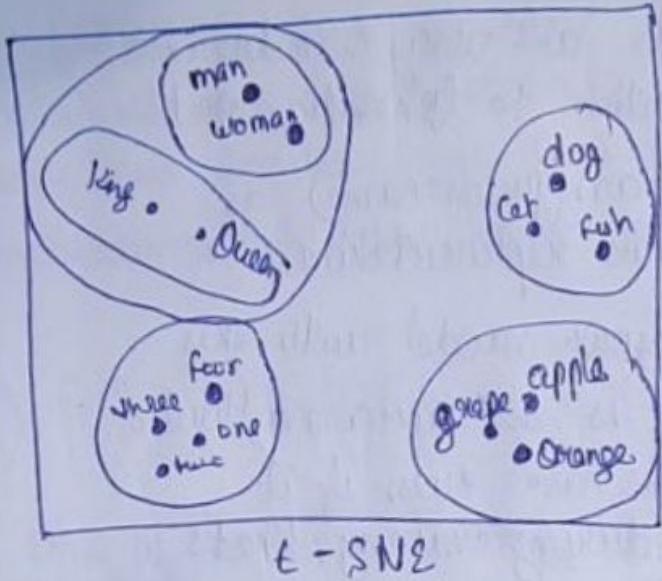
	(5391) Man	(9853) Woman	(4914) King	(9157) Queen	(4586) Apple	(657) Orange	Word Vector by itself
gender	-1	1	-0.95	0.97	0.00	0.01	
royal	0.01	0.02	0.93	0.95	-0.01	0.00	
age	0.03	0.02	0.7	0.69	0.03	-0.02	
food	0.04	0.001	0.002	0.01	0.95	0.97	
	—	—	—	—	—	—	
	—	—	—	—	—	—	
	—	—	—	—	—	—	

↑ 300 ↓ This 300 means it's 300 dimensional (having 300 features)

Now if we use this representation to represent the words orange and apple then we'll notice that the representations for orange & apple are quite similar

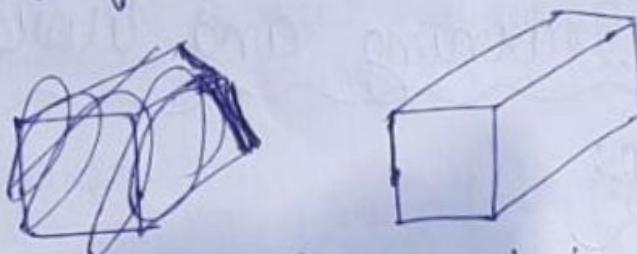
Visualizing Word embeddings

If we are able to learn a 300D feature vector or 300D embedding for each words, one of the popular things to do is to take 300D data and embed it to 2D space (so that we can visualize) one common algorithm for this is t-SNE



man and woman are grouped together
 King & Queen " " "
 so man, woman, King, Queen
 group together.

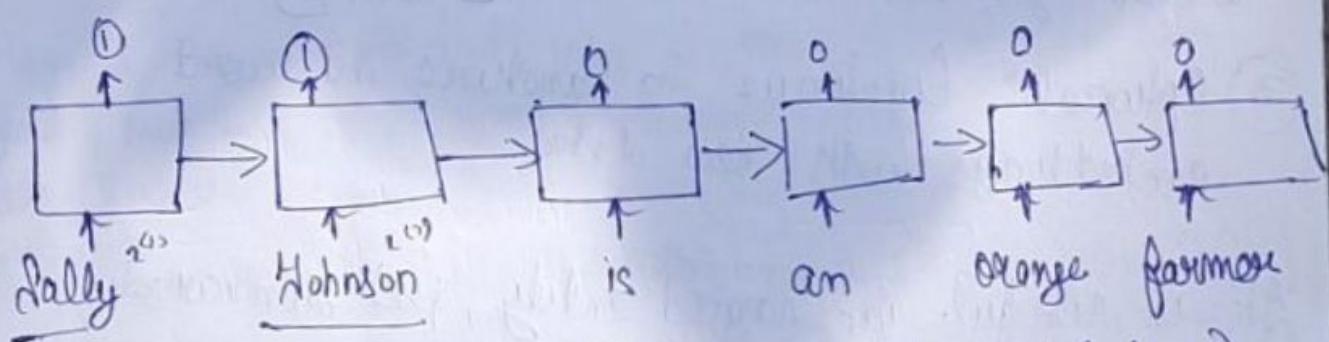
Word embedding for this example is 300D space like this



word embedding is one of the most imp. ideas in NLP

② USING WORD EMBEDDING

Linear example (for name entity recognition problem)



If we now see a new input (on our trained model alone)
 Robert I'm 'm an apple. farmer

knowing that orange and apple are very similar will make it easier for your algorithm to generalize that Robert Lin is also a human (person's name) as apple and orange have near representation.

Now if you have tested your model with this sentence "Mahmoud Badry is a durian collector" the network should learn the name even if it hasn't seen the word durian during training. That's the power of word representation

Transfer learning and Word embeddings

- ① Learn word embeddings from large datasets
(or download ~~and~~ pre-trained embeddings online)
- ② Transfer embedding to new task with smaller training set. (rather than using a 10,000 dimensional vector you can now use 300D vector)
- ③ Optional: Continue to fine-tune the word embeddings with new data.

It is useful for named entity, text summarization, etc.

In face recognition we learned about Siamese network architecture in which we encode each face into a vector and then check how similar are those vectors here encoding of embeddings have similar meaning.

③ PROPERTIES OF WORD EMBEDDINGS

One of the most fascinating properties of word embeddings is that they can also help with analogy reasoning. While analogy reasoning may not be itself the most important NLP application, but it might help convey a sense of what these word embeddings can do.

Word embedding table:

	Man (5343)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.76	0.69	-0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

$e_{\text{man}} - e_{\text{woman}}$ \approx $e_{\text{king}} - e_{\text{queen}}$

Man is to woman and King is to ____? We can find out in the word embedding table.

lets subtract $e_{\text{man}} - e_{\text{woman}}$ from e_{king} (vectors)

Well get $e_{\text{man}} - e_{\text{woman}} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -2 & 0.01 & 0.01 & 0.08 \end{bmatrix}$

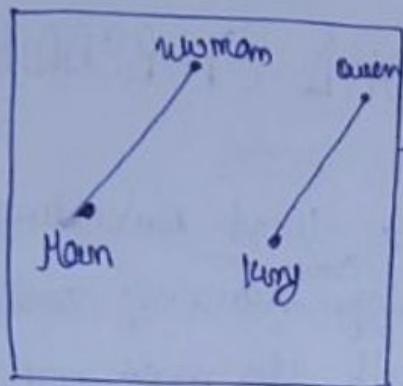
Similarly for King & Queen

$e_{\text{king}} - e_{\text{woman}} = \begin{bmatrix} -2 & 0 & 0 & 0 \end{bmatrix}$ (this is approx)

(we can also write this in column)

To find analogy we have to compare man, woman
with king & queen

eman - ewoman ≈ king - equeen (≈ this mean close^{to})



- This vector represents the gender difference
- difference is about the gender in ~~both~~ both
- ~~approximate~~
- This 2D of 4 Director using t-SNE

To find out we use this formula

• eman - ewoman ≈ ekong - equeen?

for this we have to find ^{word} w(queen) that maximizes the similarity b/w ew(woman) to ekong - eman + ewoman

$$= \underset{\text{Similarity}}{\text{Sum}}(ew, ekong - eman + ewoman)$$

It turns out that equeen is the best solution here
Solution here that gets the similar vector.

The most similarity function is called Cosine Similarity

$$\text{Sum}(U, V) = \frac{U^T V}{\|U\|_2 \|V\|_2}$$

$\boxed{U^T = \text{dot product}}$

$$U = ew$$

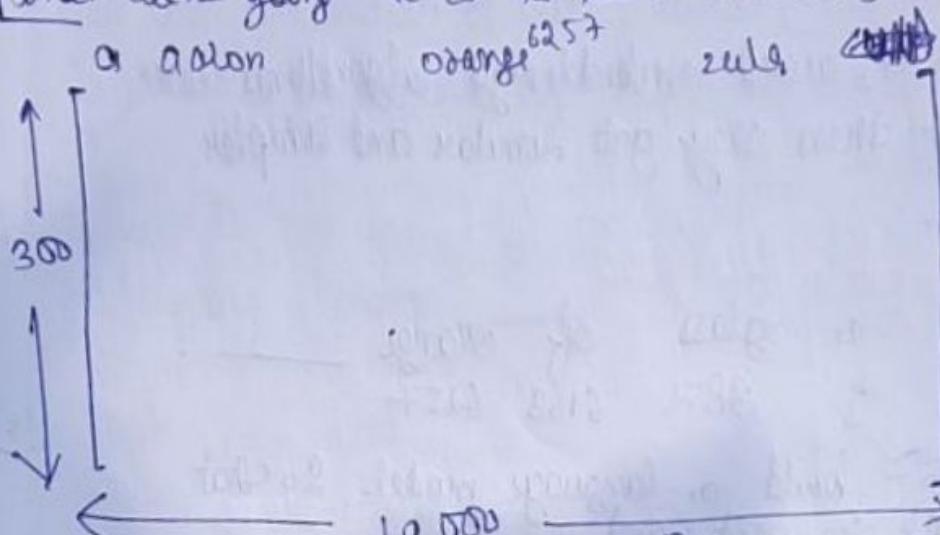
$$V = ekong - eman + ewoman$$

④ EMBEDDING MATRIX

When you implement an algorithm to learn a word embedding

Let's say we are using 10,000-word vocabulary

What we're going to do is learn embedding matrix E ,



Orange will be represented as

$$\begin{bmatrix} 0.6257 \\ \vdots \\ \vdots \\ \vdots \\ - \end{bmatrix} \rightarrow 6257 \quad \begin{array}{l} \uparrow \\ 10,000 \end{array}$$

Embed matrix

$$E \cdot O_{6257} = \begin{bmatrix} ? \end{bmatrix} \text{ given this will be a } 300 \times 1 \text{ vector}$$

$$E = (300, 10,000)$$

$$O = (10,000, 1)$$

$$(300, 1)$$

$$= e_{6257}$$

Our goal is to learn our embedding matrix E
 & we initialize E randomly and we use gradient
 descent to learn all the parameters of this 300-dimensional
 by 10,000 dimensional matrix.

To LEARNING WORD EMBEDDINGS WORD2VEC & GLOVE

① LEARNING WORD EMBEDDINGS

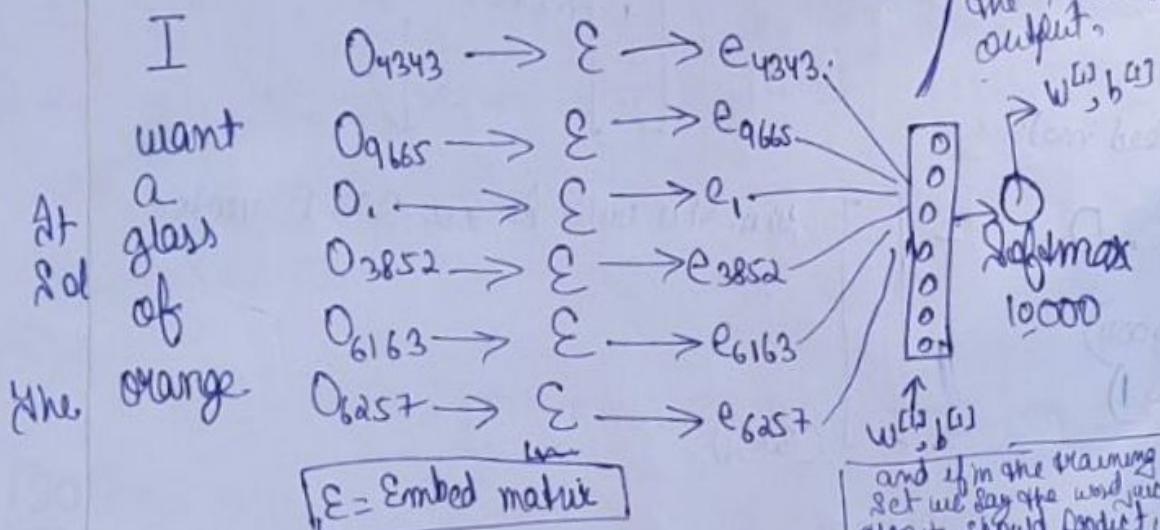
At the start, word embeddings algorithms were complex but then they got simpler and simpler.

Example:-

T₁ I want a glass of orange — .
4343 9665 1 3852 6163 6257

• we want to build a language model so that if we can predict the next word.

• So we use this neural network to learn the language model



Σ = Embed matix

and if in the training
let we say the word we
aren't should predict word just
like this

• In this example we look last 6 words that falls behind the word that we want to predict.

S! Suppose we have a example: "I want a glass of orange juice to go along with my cereal"

Now suppose we want to predict word "juice". Here we can do this by ~~many~~ ways.

a) last 4 words

We use a window of last 4 words (4 is hyperparameter)

"a glass of orange" and try to predict the next word from it.

b) 4 words on the left or right

a glass of orange _____ to go along with

c) last one word

which is "orange".

To summarize, the language modeling problem poses a machine learning problem where you input the context ~~like~~ like the last four words and predict some target words. And posing that problem allows you to learn good word embeddings.

② WORD2VEC

lets talk about skip-grams first -

for example:

our sentence \rightarrow "I want a glass of ~~my~~ orange juice to go along with my cereal".

In skip-gram we choose context and target.

In Skip gram, what we are going to do is come up with few context to target words to create our supervised learning algorithm problem. So rather than having the context be always the last four words or last ten words immediately before the target word, what we'll do is pick a word to be context word.

The target is chosen randomly based on a window with a specific size

Context	Target	How far
Orange	sue	+1
Orange	glass	-2
Orange	my	+6

we want to use this learning problem to learn word embeddings

* This is not an easy learning because learning within +10 or -10 words which is hard

Word2Vec model: \rightarrow It has 2 algorithms in it
 ① CBOW
 ② Skip gram

Let's learn basics about Word

\rightarrow Here the model learns C to t (Context to target)

* we get e_c by $E_{\theta} e_c$

$$P(t|C) = \hat{y}$$

In this algorithm use cross entropy loss function

$$P(t|C) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{V_{vocab}} e^{\theta_j^T e_c}}$$

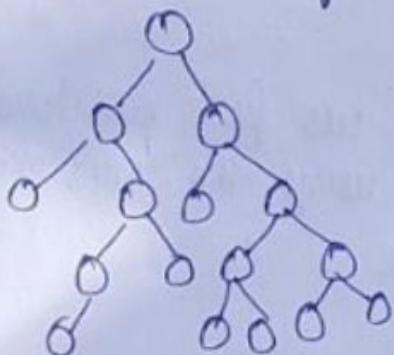
θ_t = parameter associate with output t

loss function for Softmax will be used which is

$$L(\hat{y}, y) = - \sum_{i=1}^{10000} y_i \log \hat{y}_i$$

y as a one-hot vector

- Here we are summing 10,000 numbers which correspond to the number of words in our vocabulary.
- If the no. of our vocabulary is very large say 1 million the computation will become very slow.
But we can solve the computation problem by using "Hierarchical Softmax classifier" which is tree classifier.



But in practice the hierarchical softmax classifier doesn't use a balanced tree like the drawn one. Common words are at top which are used commonly and less common are at bottom

How to sample the context?

- choose random c from corpus
- There are some words common words which dominate other words.

- At user & ideas of learning word embeddings
 - ① Skip-gram model
 - ② CBOW (continuous bag-of-words)

③ NEGATIVE SAMPLING

Negative Sampling allows you to do something similar to the skip-gram model, but with a much more efficient learning algorithm. We will create a different learning problem.

- ~~CBOW~~

Our sentence: "I want a glass of orange juice to go with my meal"

Our Sampling

Context	Word	Target
orange	juice	1
orange	king	0
orange	book	0
orange	the	0
orange	of	0

We get positive example by using the same skip-gram technique

To generate a negative example, we pick a word randomly from the vocabulary

We notice that, we got word "of" as negative example although it appeared in the same sentence

④ GLOVE WORD VECTORS

It refers to global vector for word representation.

Preciously, we were sampling pairs of words context and target words by picking 2 words that appear in close proximity to each other in our dataset corpus.

So what the glove algorithm does is, it stands off just by making that explicit.

$$X_{ij} = \# \text{times } i \text{ appears in the context of } j$$

$i \in t \quad \downarrow \quad t \quad \downarrow \quad c \in C$

$X_{ij} = X_{ji}$ if we choose a window pair, but they will not equal if we choose the previous words for for example. In glove they use

a window which means they are equal

x_{ij} is a count that captures how often do words i and j appear with each other,

So what glove model does it optimize the following

Model

minimize

APPLICATIONS USING WORD EMBEDDINGS

① SENTIMENT CLASSIFICATION

example

x

The dessert is excellent.

y

★ ○ ○ ○ ○

Service was quite slow.

○ ○ ○ ○ ○

One of the problem with the sentiment classification is you might not have a huge labeled data set.

RNN is used for solving problem of Sentimental Analysis

Sentence: "Completely lacking in good taste, good service, and good ambience"

Here if we do a simple classification for this sentence it won't work. So, we use RNN for this

② DEBIASING WORD EMBEDDING

In this we'll learn the ideas for diminishing or eliminating these forms of bias in word embeddings. Here, bias doesn't mean one bias variants. It means gender, ethnicity, sexual orientation bias.

The problem of bias in word embeddings:

- Man:Woman as King:Queen ✓
- Man:Computer Programmer as Woman: Housemaker X
- Father: Doctor as Mother:Nurse X

Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the data used to train the model.

Addressing Bias in word embeddings

So here's the idea. Let's say we've already learned a word embedding like this -



Given above learned word embeddings:

• doctor

babysitter

grandmother

• grandfather

girl •

• Body

• he

she •

M Do, first thing we are going to do it is identify the direction corresponding to a particular bias we want C to reduce or eliminate

C Do how do you identify the direction

X Do how do you identify the direction for the case of gender what we can do is take the embedding vector for d he and subtract the embedding vector for d she, because the differs by gender

W Here are the steps:

a.) Identify the direction:

o male - female

o female - female

o =

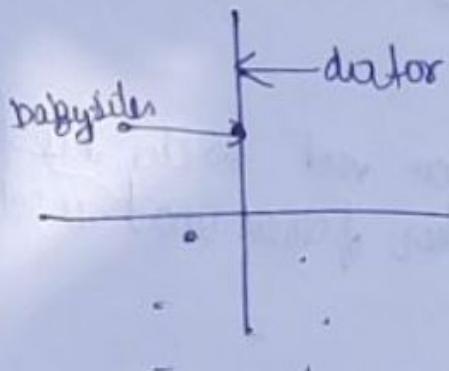
o choose some K differences & average them

o this will help you to find this

babysitter		doctor
grandmother	grandfather	bias
girl		• boy
he		• she

By that we found the bias direction which is Doctor and the non-bias vector which is girl vector.

- b) Neutralize: for every word that not defintional, project do get rid of bias
- Babysitter & doctor need to be neutralized so we project them on non-bias axis with direction of the bias

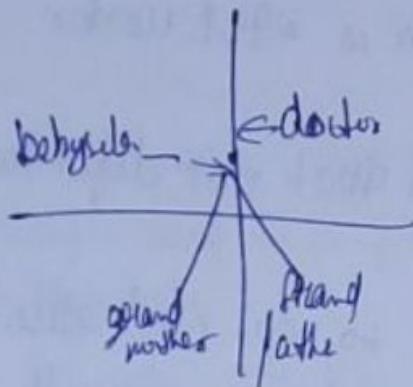


after that they will be equal in the term of gender. To do this the authors of the paper trained a classifier to tell the words that needs to be neutralized or not.

of Equalise. Part

- ☒ we want each pair to have difference only in gender like
 - Grandfather - Grandmother He - She
- ☒ we want to do this because the distance b/w Grandfather & babykiller is bigger than babykiller & grandmother

(c)



- ☒ To do that we move grandfather & grandmother to a point where they will be in the middle of the no-bias axis
- ☒ There are some words you need to do this for in your steps. No Number of other word is related. Small