

DEEP LEARNING - course 1

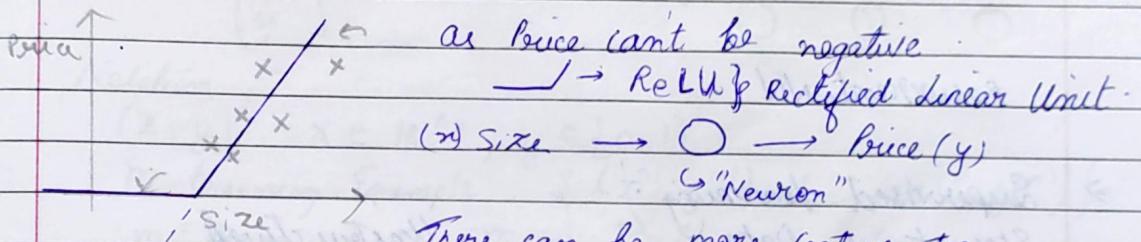
Page No. 101
Date 9.5.21

→ AI is the new electricity and will bring about an equally big transformation.

* What we'll learn in this specialization?

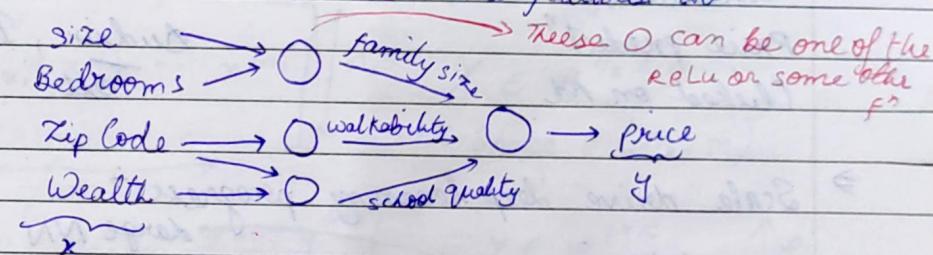
1. Neural Network & Deep Learning
2. Improving Deep Neural Networks - Hyperparameter tuning, Regularization and optimization
3. Structure your Machine Learning project
4. Convolutional Neural Networks [CNN]
5. Natural language processing - Building sequence models (RNN, LSTM)

⇒ House Price Prediction :-



(*) Size → ○ → Price (y)
"Neuron"

There can be more features too

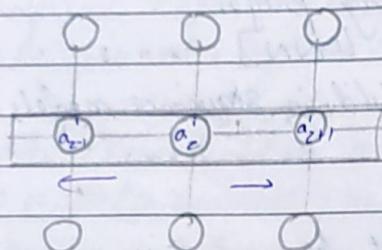
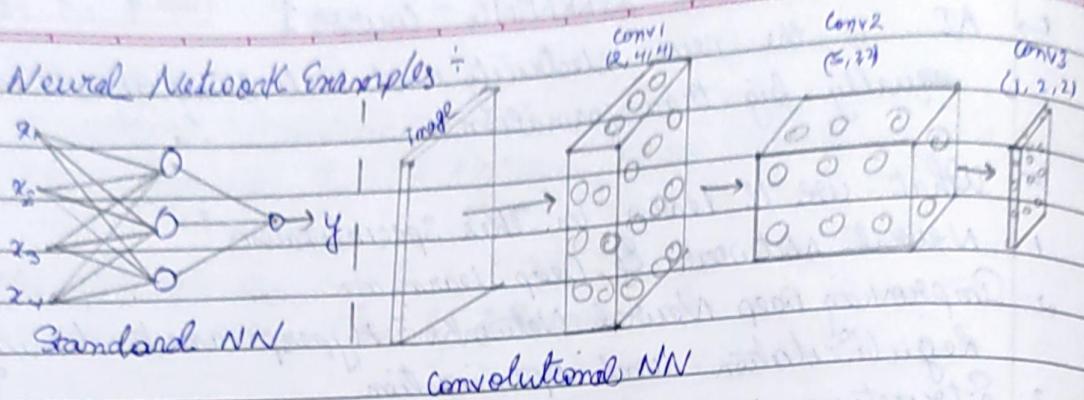


⇒ Supervised Learning :-

Input (x)	Output (y)	Application
Home features	Price	Real Estate
Ad, user info	Click on Ad (0,1)	Online Advertising
Image	Object (1, ..., 1000)	Photo Tagging
Audio	Text transcript	Speech Recognition
English	Chinese	Machine Translation
Image, Radar info	Position of other cars	Autonomous Driving

{ standard NN CNN RNN custom / hybrid }

⇒ Neural Network Examples :-



Recurrent NN

⇒ Supervised learning :-

Structured Data

Price predict

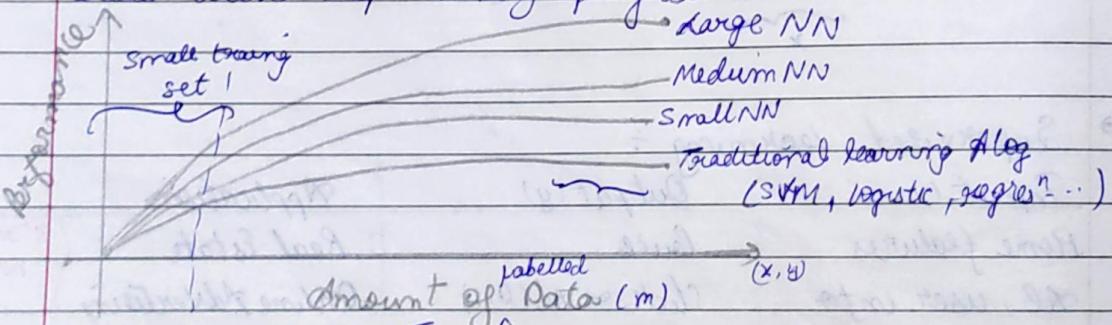
Clicked on Ad

Unstructured

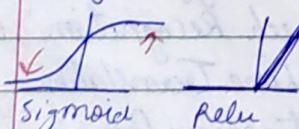
Audio, Image,

Text

⇒ Scale drive deep learning progress :-



Algorithms :- To run faster ReLU is better than Sigmoid?



In Sigmoid slope is close to zero at end so

learning becomes very slow. When GD is used the parameters change very slowly; when we use ReLU gradient = 1, so, GD runs faster than sigmoid. Computationally cheap

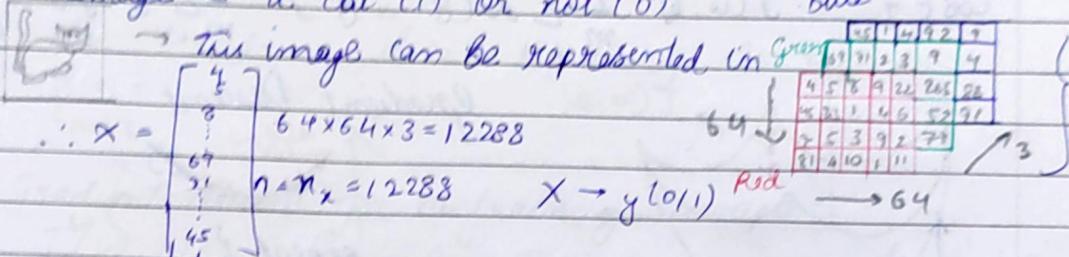
The cycle can be done in seconds, minutes, days or even weeks. \rightarrow Idea, Experiment, Code.

Week - 2

Basics of Neural Network

Binary Classification

Ex. Image is a cat (1) or not (0) \rightarrow Blue



Notation :-

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

m training Example: $\{ (x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)}) \}$

$m = M_{train}$; $M_{test} \Rightarrow$ test examples

$$X = \begin{bmatrix} | & & & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} & | \\ | & & & & | \\ \vdots & \vdots & & \vdots & \vdots \\ m & & & & \end{bmatrix} \quad X \in \mathbb{R}^{n_x \times m}$$

$$Y = [y^{(1)} \ y^{(2)} \ \cdots \ y^{(m)}] \quad ; \quad Y \in \mathbb{R}^{1 \times m}$$

Y shape = $(1, m)$

\Rightarrow Logistic Regression :-

Given X ; we want $\hat{y} = P(y=1|x)$, $x \in \mathbb{R}^{n_x}$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

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

If z is large (≈ 1) we do use $x_0 = 1$, $x \in \mathbb{R}^{n_x+1}$

$$\text{If } z \text{ is small } (\approx 0), \hat{y} = \sigma(\theta^T x) \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \quad w$$

⇒ Cost function :-

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

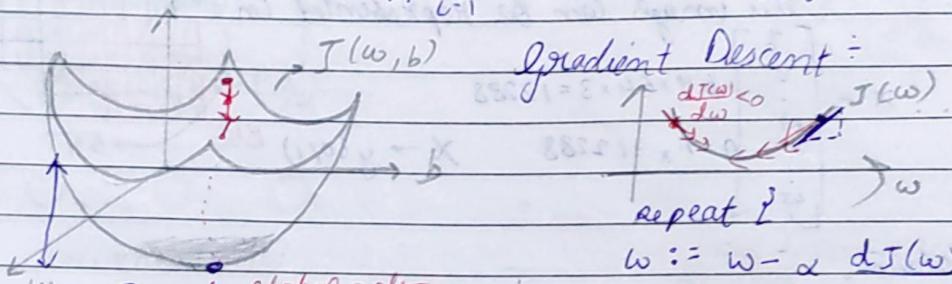
Applied to single training example

$$L(\hat{y}, y) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})] \quad \{ \text{cost f^n}$$

If $y=1$ $L(\hat{y}, y) = -\log \hat{y} \leftarrow$ as small as possible i.e. yes

If $y=0$ $L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow$ y must be big

$$\text{Cost f^n : } J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})$$



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

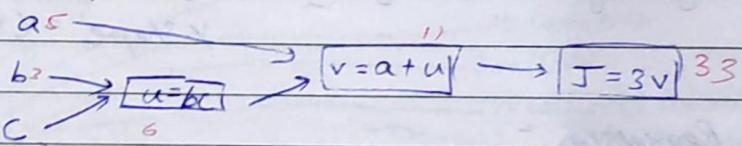
For graph purpose we are taking $b = b - \alpha \frac{dJ(w, b)}{db}$
it as a single vector, here

⇒ Computational graph :-

$$\text{Say } J(a, b, c) = 3(a + bc) \quad u = bc \quad a = 5$$

$$\text{Now how to compute } v = au \quad b = 3$$

$$\text{using 3 steps} \quad \underbrace{v}_{\substack{v \\ \text{ }} \rightarrow \text{ }} \quad \underbrace{J}_{\substack{J \\ \text{ }} \rightarrow \text{ }} \quad J = 3v \quad c = 2$$



$$\text{lets find } \frac{dJ}{dv} = ? \quad ; \quad \frac{dJ}{da} = ?$$

⇒ Logistic Regression :-

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

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

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

$$L(a, y) \quad \{ \frac{\partial L(a, y)}{\partial a} \leftarrow \frac{y - \hat{y}}{a(1-a)}$$

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

$$w_1 \rightarrow \uparrow \frac{\partial L(a, y)}{\partial z}$$

$$x_2 \rightarrow \uparrow \frac{\partial L(a, y)}{\partial z}$$

$$w_2 \rightarrow \uparrow \frac{\partial L(a, y)}{\partial z}$$

$$\hat{y} = a = \sigma(z) \quad \uparrow \frac{\partial L(a, y)}{\partial a}$$

$$\frac{-y}{a} + \frac{1-y}{1-a}$$

$$\frac{dL}{dz} = \frac{dL}{da} \cdot \frac{da}{dz} \rightarrow a(1-a) \rightarrow \text{sigmoid derivative}$$

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

$$\frac{dL}{dz} = \frac{-y + ay + 1 - ay}{a(1-a)} = \frac{1 - y}{a(1-a)}$$

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

These are
Backpropagation
step
to go back

$$\frac{dL}{d\omega_1} = x_1 dz; \quad d\omega_2 = x_2 dz; \quad db = dz$$

\Rightarrow Gradient descent for m training example:

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)}) \quad d\omega_1^{(i)}, d\omega_2^{(i)}, db^{(i)}$$

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

$$\therefore \frac{\partial J(\omega, b)}{\partial \omega_1} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L(a^{(i)}, y^{(i)})}{\partial \omega_1}$$

$$\underbrace{d\omega_1^{(i)}}_{\rightarrow (x^{(i)}, y^{(i)})} \rightarrow (x^{(i)}, y^{(i)})$$

\Rightarrow Logistic Regression on m Examples:

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

First for loop
For $i = 1$ to m :

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

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

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

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

$$d\omega_1 += x_1^{(i)} \cdot dx^{(i)} \quad \} \text{for } n=2$$

$$d\omega_2 += x_2^{(i)} \cdot dx^{(i)}$$

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

$$J/ = m; \quad d\omega_1 / m; \quad d\omega_2 / = m; \quad db / = m;$$

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

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

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

Second for loop
though we only
choose 2 features

★ Vectorization :-

What is Vectorization ?

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

Non-vectorized

$$z = 0$$

for i in range(n_x):

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

$$z += b$$

$$w = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, x = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, z, w \in \mathbb{R}^{n_x}$$

Vectorized

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

$$\underbrace{w^T x}$$

* How much fast vectorized version is

Create your own Jupyter Notebook

CPU, GPU → SIMD - single instruction multiple data

⇒ Neural Network programming guidelines :-

Whenever possible, avoid explicit for-loops.

$$u = A v$$

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

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

for i ...

for j ...

$$u[i] += A[i][j] * v[j]$$

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

*Vectorized
way faster*

Ex

We need exponential operation on each and every element

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

For-loops

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

for i in range(n):

$$u[i] = \text{math.exp}(v[i])$$

Vectorized

import numpy as np

$$u = np.exp(v)$$

⇒ Logistic regression derivatives :

$$J=0 \quad (\text{d}w_1=0, \text{d}w_2=0), \text{d}b=0 \quad \text{d}w = np \cdot \text{softmax}(b_{x,1})$$

are for
loop

for $i \in \text{range}(m)$:

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

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

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

$$\text{Second for loop } d z^{(i)} = a^{(i)}(1-a^{(i)})$$

~~$$\text{is REMOVED & vectorized method is taken}$$~~

$$d w_1 += x_1^{(i)} \cdot dz^{(i)} \quad d w += x^{(i)} dz^{(i)}$$

~~$$d w_2 += x_2^{(i)} \cdot dz^{(i)}$$~~

~~$$d b += dz$$~~

$$J = J/m; \quad \boxed{d w_1 = d w_1 / m; \quad d w_2 = d w_2 / m; \quad d b = d b / m}$$

$$d w / m$$

⇒ Vectorizing Logistic Regression =

$$z^{(1)} = w^T x^{(1)} + b \quad z^{(2)} = w^T x^{(2)} + b \quad z^{(3)} = w^T x^{(3)} + b$$

$$a^{(1)} = \sigma(z^{(1)}) \quad a^{(2)} = \sigma(z^{(2)}) \quad a^{(3)} = \sigma(z^{(3)})$$

$$X = \begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & 1 \end{bmatrix}; \quad X \in \mathbb{R}^{n_x \times m} \quad w^T \begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & 1 \end{bmatrix}$$

$$Z = [z^{(1)} \ z^{(2)} \dots z^{(m)}] = w^T X + [b \ b \ \dots \ b] \quad 1 \times m$$

$$= [w^T x^{(1)} + b \ w^T x^{(2)} + b \ \dots \ w^T x^{(m)} + b] \rightarrow 1 \times m$$

$$Z = np \cdot \text{dot}(w^T, X) + b \rightarrow (1, 1)$$

$\underbrace{1 \times m}_{\infty}$ "Broadcasting"

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

So here we didn't use any for loops for our model.

⇒ Vectorizing Logistic Regression's Gradient Computation =

$$d z^{(1)} = a^{(1)} - y^{(1)}; \quad d z^{(2)} = a^{(2)} - y^{(2)} \dots$$

$$d Z = [d z^{(1)} \ d z^{(2)} \ \dots \ d z^{(m)}] \quad 1 \times m$$

$$A = [a^{(1)} \ \dots \ a^{(m)}] \quad d Z = A - Y = \underline{\underline{[a^{(1)} - y^{(1)} \ \dots \ a^{(m)} - y^{(m)}]}}$$

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

we still have this for loop left:

$$\left. \begin{array}{l} dw = 0 \\ dw += x^{(i)} dz^{(1)} \\ dw += x^{(i)} dz^{(2)} \\ dw / = m \end{array} \right\} \quad \left. \begin{array}{l} db = 0 \\ db += dz^{(1)} \\ db += dz^{(2)} \\ db += dz^{(m)} \\ db / = m \end{array} \right\}$$

In one line:

$$\left. \begin{array}{l} db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} = np \cdot \text{sum}(dz) \end{array} \right\}$$

$$\left. \begin{array}{l} dw = \frac{1}{m} X \cdot dz^T = \frac{1}{m} \left[\begin{array}{c} x^{(1)} \\ \vdots \\ x^{(m)} \end{array} \right] \cdot \left[\begin{array}{c} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{array} \right] \\ = \frac{1}{m} [x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)}] \rightarrow n \times 1 \end{array} \right\}$$

\Rightarrow Vectorized Implementation for iter in range(1000):

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

We have $A = \sigma(z)$

to use $dz = A - y$

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

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

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

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

Can't replace it

\Rightarrow Broadcasting in python:

Ex Calories from Carbs, Proteins, Fats in 100g of different food

Apples Beef Eggs Potatoes

Carbs [56.0 0.0 4.4 68.0]

Protein [1.2 104.0 52.0 0.8]

Fat [1.8 135.0 99.0 0.9]

calculate % of carbs from

A Carb, Protein, Fat!

(3, 4)

$\overbrace{\text{59 cal}}$

$\frac{56}{59} \approx 95\% \text{ of carb}$

Let's do it without a for-loop?

`call = A.sum(axis=0)` → for column axis=0

`percentage = 100 * A / (call.reshape(1, 4))` • not necessary to call

* check for Broadcasting in Jupyter Notebook.

$(3, 4) / (1, 4)$ matrix.

$$\text{Ex: } \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 \xrightarrow{\text{Python Broadcast}} \begin{bmatrix} 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

$(1, 1) \text{ to } (3, 1)$ Matrix

$$\text{Ex: } \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$(m, n) (2, 3)$ Python Broadcasted $(1, 3)$ to $(2, 3)$ Matrix

$$\text{Ex: } \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

$(2, 3) (2, 1) \rightarrow (2, 3)$

⇒ General principle :

$$(m, n) \xrightarrow{\text{Matrix}} (1, n) \xrightarrow{\text{Broadcast to}} (m, n)$$

$$(m, 1) \xrightarrow{\text{Matrix}} (m, 1) \xrightarrow{\text{Broadcast to}} (m, n)$$

$$\text{Ex: } (m, 1) \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

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

Matlab/Octave = `bzfun` {Broadcasting in Matlab/octave}

⇒ Python/numpy vectors :

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

X `a.shape = (5,)` → This is a "rank 1 array" { Don't use rank 1 arrays }

✓ `a = np.random.rand(5, 1)` { use this column vector }

`a.shape = (5, 1)`

`assert(a.shape == (5, 1))` { To check for shape }

`a = a.reshape(5, 1)` { You can reshape rank 1 array to column vector }

* Vectors/Matrix in Python & Reshaping . . .

⇒ Logistic Regression Cost function :

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

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

$$\text{If } y=1 : P(y|x) = \hat{y} \quad \left. \right\} P(y|x)$$

$$\text{If } y=0 : P(y|x) = 1 - \hat{y} \quad \left. \right\} P(y|x)$$

$$\therefore P(y|x) = \hat{y}^y \cdot (1-\hat{y})^{(1-y)}$$

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

$$\text{If } y=0 \Rightarrow P(y|x) = (1-\hat{y})$$

$$\begin{aligned} \log P(y|x) &= \log \hat{y}^y \cdot (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ &= -d(\hat{y}, y) \end{aligned} \quad \text{minimum cost f^n}$$

⇒ Cost on m examples :

$$P(\text{labels in training set}) = \prod_{i=1}^m P(y^{(i)}|x^{(i)})$$

$$\log P(\dots) = \sum_{i=1}^m \underbrace{\log P(y^{(i)}|x^{(i)})}_{-d(\hat{y}^{(i)}, y^{(i)})} \quad \left. \begin{array}{l} \text{Maximum Likelihood} \\ \text{Estimation} \end{array} \right.$$

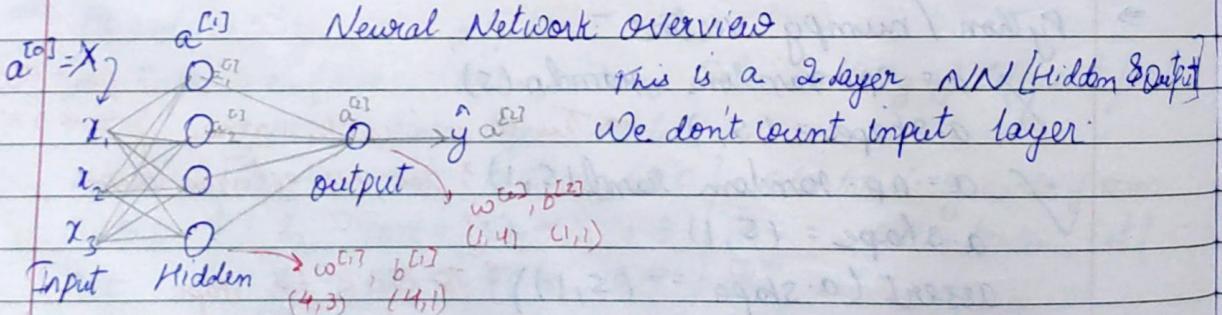
$$= - \sum_{i=1}^m d(\hat{y}^{(i)}, y^{(i)})$$

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

(Minimize)

Week - 3

Neural Network overview



$x_1, x_2, x_3 \rightarrow \hat{y}$

$$\text{2 steps take place: } z^{[1]} = w^{[1]T} x + b^{[1]} \quad | \quad \begin{array}{l} z^{[1]} = w^{[1]T} x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \end{array} \quad | \quad \begin{array}{l} \text{For first hidden node} \\ a^{[1]} \in \text{layer} \\ a^{[1]} \in \text{node in layer} \end{array}$$

$$z^{[1]} = w^{[1]T} x + b^{[1]} ; a^{[1]} = \sigma(z^{[1]}) \quad | \quad \begin{array}{l} \text{Instead of using for-loop} \\ \text{We will be vectorising these equations.} \end{array}$$

$$z^{[2]} = w^{[2]T} x + b^{[2]} ; a^{[2]} = \sigma(z^{[2]})$$

$$z^{[3]} = w^{[3]T} x + b^{[3]} ; a^{[3]} = \sigma(z^{[3]})$$

$$z^{[4]} = w^{[4]T} x + b^{[4]} ; a^{[4]} = \sigma(z^{[4]})$$

$$z = \begin{bmatrix} z^{[1]} \\ z^{[2]} \\ z^{[3]} \\ z^{[4]} \end{bmatrix} = \underbrace{\begin{bmatrix} \dots & w_1^{[1]T} \\ \dots & w_2^{[1]T} \\ \vdots & \vdots \\ \dots & w_4^{[1]T} \end{bmatrix}}_{w^{[1]}} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ \vdots \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ \vdots \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix}$$

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

$$= \underbrace{w^{[1]}}_{(4 \times 3)}$$

$$b^{[1]}$$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]}) \quad | \quad a^{[0]} \quad | \quad a^{[1]} = \sigma(z^{[1]})$$

$$\therefore z^{[1]} = w^{[1]T} x + b^{[1]} \quad | \quad (4,1) \quad | \quad (4,1)$$

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

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

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

\Rightarrow Vectorizing across multiple examples:

$$x \rightarrow a^{[1]} = \hat{y}$$

for various x cases

$$x^{(1)} \rightarrow \hat{y}^{(1)} a^{[2](1)} \quad | \quad a^{[2](1)} \rightarrow \text{example } i$$

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

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

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

for $i = 1 \dots m$,

$$z^{[1](i)} = w^{[1]} x^{(i)} + b^{[1]} \quad | \quad z^{[1]} = \begin{bmatrix} z^{1} \\ \vdots \\ z^{[1](m)} \end{bmatrix}$$

$$a^{[1](i)} = \sigma(z^{[1](i)}) \quad | \quad \begin{array}{l} \text{diff features} \\ \text{feature} \end{array}$$

$$z^{[2](i)} = w^{[2]} a^{[1](i)} + b^{[2]} \quad | \quad A^{[2]} = \begin{bmatrix} a^{[2](1)} \\ \vdots \\ a^{[2](m)} \end{bmatrix}$$

$$a^{[2](i)} = \sigma(z^{[2](i)}) \quad | \quad \begin{array}{l} \text{Hidden unit-number} \\ \text{!} \end{array}$$

$\xrightarrow{\text{training example}}$

So our Vectorized Algo will be:

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

$$A^{[0]} = \sigma(Z^{[0]})$$

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

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

\Rightarrow Justification for vectorized implementation:

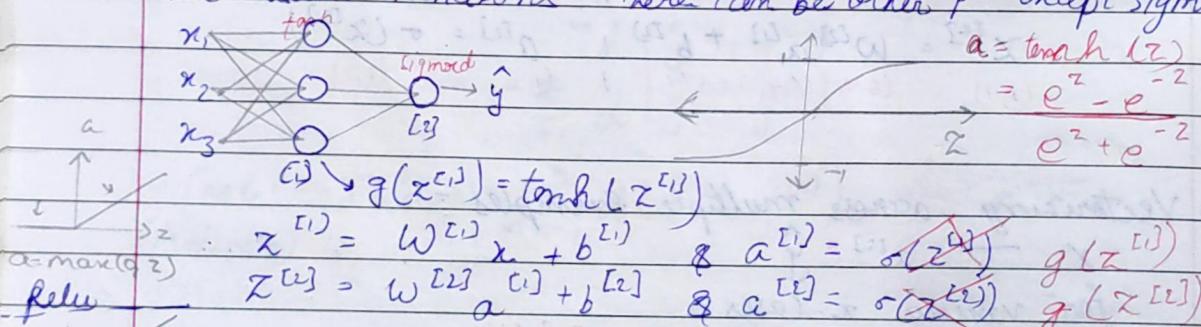
$$Z^{[1](i)} = w^{[1]T} X^{(i)} + b^{[1]} ; Z^{[0](j)} = w^{[0]T} X^{(j)} + b^{[0]}$$

$$w^{[1]} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} ; w^{[1]} X^{(i)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} ; w^{[0]} X^{(j)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$Z^{[1]} = w^{[1]} \begin{bmatrix} x^{(i)} & x^{(0)} & \dots \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix} = \begin{bmatrix} Z^{[1](i)} \\ Z^{[1](0)} \\ \dots \end{bmatrix}$$

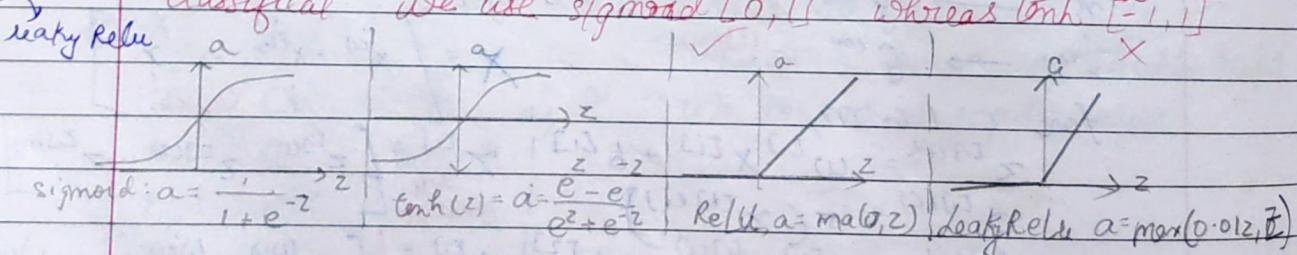
(We can add $b^{[1]}$ via broadcasting $+ b^{[1]} + b^{[0]} \dots$)

\Rightarrow Activation Functions: There can be other f except sigmoid.



$$\text{relu} \quad z^{[1]} = w^{[1]} x + b^{[1]} \quad \& \quad a^{[1]} = \sigma(z)$$

We mostly use tanh; just for final node of binary classification we use sigmoid [0, 1] whereas tanh [-1, 1]



⇒ Why do we need Non-Linear Activation function?

Why we are using f^n like sigmoid, tanh, or ReLU

We know $a^{[L]} = g^{[L]}(z^{[L]})$ instead we will write

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

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

$$\text{or } a^{[2]} = w^{[2]} + w^{[2]}a^{[1]} + b^{[2]}$$

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

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

If we use linear f^n we are eventually getting a linear f^n as our output. This functⁿ is only useful when we predict house price. $g(z) = z$ and that too at the output layer not in between. We can also use ReLU for housing price prediction.

⇒ Derivatives of Activation functⁿ:

i) Sigmoid f^n $f(z) = \frac{1}{1+e^{-z}}$ $f'(z) = \frac{0+(-e^{-z})}{(1+e^{-z})^2}$

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

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

$$g'(z) = a(1-a)$$

ii) $g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ $\therefore g'(z) = \frac{(e^z + e^{-z})(e^z - e^{-z}) - (e^z - e^{-z})^2}{(e^z + e^{-z})^2}$

$$g'(z) = \frac{2e^z}{e^z + e^{-z}} = \frac{e^z + e^{-z} + e^z - e^{-z}}{e^z + e^{-z}} = 1 + \frac{(e^z - e^{-z})^2}{(e^z + e^{-z})^2}$$

$$g'(z) = 1 - (\tanh(z))^2 \quad \therefore \text{If } z=0; \tanh \approx 1; g'(z) \approx 0$$

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

$$\underline{\underline{g'(z) = 1 - a^2}}$$

$$z=0; \tanh \approx 0; g'(z)=1$$

iii) ReLu & Leaky ReLu :

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

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

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

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

Gradient Descent for neural Network :-

$$\text{Parameters} = \omega^{(1)}, b^{(1)}, \omega^{(2)}, b^{(2)}; n_x = n^{(0)}, n^{(1)}, n^{(2)} =$$

$$(n^{(1)}, 1), (n^{(2)}, 1), (n^{(3)}, 1)$$

$$\text{cost function : } J(\omega^{[1]}, b^{[1]}, \omega^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m d(\hat{y}_i, y_i)$$

Gradient descent : Repeat {

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

$$dw^{[1]} = \frac{\partial J}{\partial w^{[0]}} ; db^{[1]} = \frac{\partial J}{\partial b^{[1]}}$$

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

$$b^{[1]} = b^{[4]} - \alpha \delta b^{[2,1]}$$

; far $\omega^{[2]}$ & $b^{[2]}$ {

⇒ Formulas for computing derivatives :-

Forward Propagation

$$\hookrightarrow z^{[1,3]} = 112^{[1,3]} x + 6^{[1,3]}$$

$$\hookrightarrow A^{\{3\}} = g^{\{3\}} / z^{\{$$

$$\Leftrightarrow Z^{[2]} = (1)^{[2]} + b^{[2]}$$

$$\hat{A}^{[2]} = g^{[2]}(z^{[2]}) = \sigma(z^{[2]}) \rightarrow db^{[2]} = \frac{1}{m} np \cdot \text{sum}(dz^{[2]}, axis=1, keepdims=True)$$

Backpropagation Step

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

$$\hookrightarrow d\omega^{[2]} = \underbrace{\star}_{m} dz^{[2]} A^{[1]} T$$

1 2 3

→ avoids Hank / array

$$\int f dz^{[k]} = \underbrace{w^{[k]T} dz^{[k]}}_{(n^{[k]}, m)} * g^{[k],'}(z^{[k]})$$

↳ element wise \times

$$\nabla \omega^{[j]} = \frac{1}{m} d z^{[j]} x^T$$

$\text{d}b^{[0,1]} = \frac{1}{m} \text{np.sum}(\text{d}z^{[1]}, \text{axis}=1, \text{keepdims=True})$
 $(0^{[1]}, 1)$ or we can keep reshape too

⇒ Backpropagation Intuition :

Computing Gradient
Logistic Regression

 x

$$w \rightarrow z = w^T x + b \rightarrow a = \sigma(z) \rightarrow d(a, y)$$

$$b \quad dz = a - y \quad da = \frac{d}{da} [d(a, y) = -y \log a - (1-y) \log(1-a)]$$

$$dw = dz \cdot X$$

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

$$db = dz$$

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

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \\ \frac{\partial a}{\partial z} = g(z) = g'(z) \end{array} \right.$$

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

$$\tilde{da}$$

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

⇒ Neural Network gradient :

$$w^{[2]} \quad \downarrow \quad \frac{d w^{[2]}}{d b^{[2]}}$$

$$x \quad \begin{aligned} dw^{[2]} &\rightarrow z^{[2]} = w^{[2]} X + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow d(a^{[2]}, y) \\ db^{[2]} &\rightarrow dz^{[2]} \quad \frac{d z^{[2]}}{d a^{[1]}} \quad da^{[1]} \quad d z^{[2]} = a^{[2]} - y \quad d a^{[1]} \end{aligned}$$

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

$$\begin{aligned} \frac{d w^{[2]}}{d b^{[2]}} &= dz^{[2]} \cdot a^{[1] T} \rightarrow dw = dz \cdot X \\ dz^{[2]} &= dz^{[2]} \end{aligned}$$

Transpose

only

$$x_1 \rightarrow 0 \rightarrow \vdots \quad \therefore w^{[2]} = (n^{[2]}, n^{[1]})$$

$$x_2 \rightarrow 0 \rightarrow \vdots \rightarrow y \quad z^{[2]}, dz^{[2]} = (n^{[2]}, 1) \text{ i.e. } (1, 1) \text{ here}$$

$$x_3 \rightarrow 0 \rightarrow n^{[2]} = 1 \quad z^{[2]}, dz^{[2]} = (n^{[2]}, 1)$$

$$x = n^{[0]} \quad n^{[0]} \quad dz^{[2]} = w^{[2] T} dz^{[1]} * g^{[2]'}(z^{[2]})$$

$$\text{similar} \quad (n^{[2]}, 1) \quad (n^{[0]}, n^{[1]}), (0^{[2]}, 1) \quad (n^{[2]}, 1)$$

$$\begin{cases} dw^{[1]} = dz^{[1]} \cdot X^T \\ db^{[1]} = dz^{[1]} \end{cases}$$

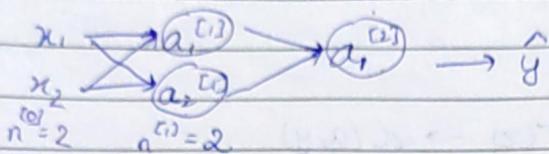
⇒ Vectorized Implementation :

All the same formula

These 6 are the vectorized formulas

$\frac{1}{m}$ appears just because $J(\hat{y}, y) = \frac{1}{m} \sum$

⇒ Random initialization : What happens if you initialize weights to zero ?



$\{ w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \}$ We have initialized it with zeros.

$$a_1^{[1]} = a_2^{[1]} ; \quad d_z^{[1]} = d z_2^{[1]}$$

$$d_w = \begin{bmatrix} 0 & 0 \end{bmatrix} ; \quad w^{[1]} = w^{[1]} - \alpha d_w \quad \& \quad w^{[1]} = \begin{bmatrix} \dots & \dots \end{bmatrix}$$

same

This will make all nodes of hidden layer identical & they will compute same thing again & again.

So we must randomly initialize the weights.

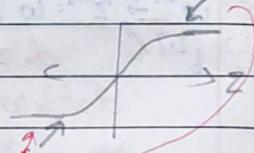
$w^{[1]} = np.random * randn(2, 2) * 0.01$ values. To have small

$b^{[1]} = np.zeros(2, 1)$ No problem with b it can be 0

Same for $w^{[2]}$ & $b^{[2]}$

If w is too large $z^{[1]} = (w^{[1]})x + b^{[1]}$

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



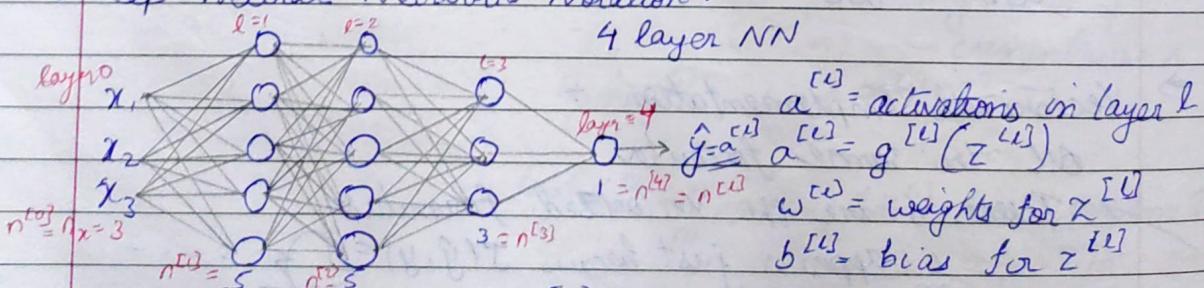
then activation unit have to start from far away so that points which will slow down

gradient descent & we keep small weights to increase speed of GD.

Week - 4

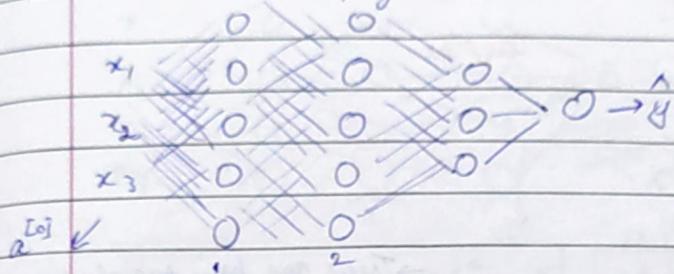
* Deep L-layer Neural Network

⇒ Deep Neural Network Notation :



$l = 4$ [no. of layers], $n^{[l]}$ = no. of units in layer l .

→ Forward propagation on a deep network =



$$x^{[l]} = w^{[l]} X + b^{[l]}$$

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

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

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

$$x^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}$$

$$a^{[L]} = g^{[L]}(x^{[L]}) = \hat{y}$$

$$\begin{aligned} x^{[l]} &= w^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} &= g^{[l]}(x^{[l]}) \end{aligned} \quad \left. \begin{array}{l} \text{General Term} \\ \hline \end{array} \right.$$

→ Vectorized =

$$\begin{bmatrix} x^{[0]} \\ \vdots \\ x^{[L-1]} \end{bmatrix} = Z^{[L]} = w^{[L]} X + b^{[L]} ; A^{[0]} = X \quad \left| \begin{array}{l} \text{for } l=1 \dots L \\ \downarrow \\ x^{[l]} = g^{[l]}(Z^{[l]}) \\ \hat{y} = g^{[L]}(Z^{[L]}) = A^{[L]} \end{array} \right.$$

→ getting your matrix dimensions right =

parameters $w^{[l]} \& b^{[l]}$

$$\begin{array}{c} \text{Dimensions: } \\ \text{Input } X \text{ is } (n^{[0]}, n^{[1]}) \\ \text{Hidden layer } l=1 \text{ is } (n^{[1]}, n^{[2]}) \\ \text{Hidden layer } l=2 \text{ is } (n^{[2]}, n^{[3]}) \\ \text{Output layer } l=L \text{ is } (n^{[L]}, 1) \end{array}$$

$$w^{[2]} = (4, 5)$$

$$w^{[3]} = (2, 4) \quad \left. \begin{array}{l} \text{dim } w^{[l]} = (n^{[l]}, n^{[l+1]}) \\ \text{dim } w^{[0]} = (n^{[0]}, n^{[1]}) \end{array} \right.$$

$$w^{[4]} = (1, 2) \quad \left. \begin{array}{l} \text{dim } w^{[l]} = (n^{[l]}, n^{[l+1]}) \\ \text{dim } w^{[0]} = (n^{[0]}, 1) \end{array} \right.$$

⇒ When using "back-propagation" $\dim w = \dim(w) = (n^{[L]}, n^{[L-1]})$

⇒ When using "back-propagation" $\dim b = \dim(b) = (n^{[L]}, 1)$

→ Vectorized = $Z^{[l]} = w^{[l]} X + b^{[l]}$

$$\begin{bmatrix} Z^{[0]} \\ \vdots \\ Z^{[L-1]} \end{bmatrix} \leftarrow \begin{array}{c} (n^{[0]}, m) \quad (n^{[0]}, n^{[1]}) \\ \times (n^{[1]}, m) \quad \rightarrow (n^{[1]}, 1) \end{array} \xrightarrow{\text{forward}} (n^{[L]}, m)$$

$$Z^{[L]}, A^{[L]} = (n^{[L]}, m) ; \text{ if } l=0 \quad A^{[0]} = X = [n^{[0]}, m]$$

$$dZ^{[L]}, dA^{[L]} = (n^{[L]}, m)$$

⇒ Why deep - neural Representations :-



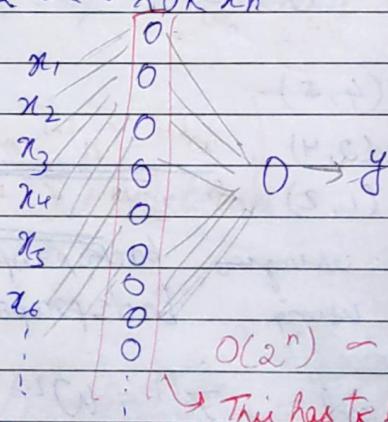
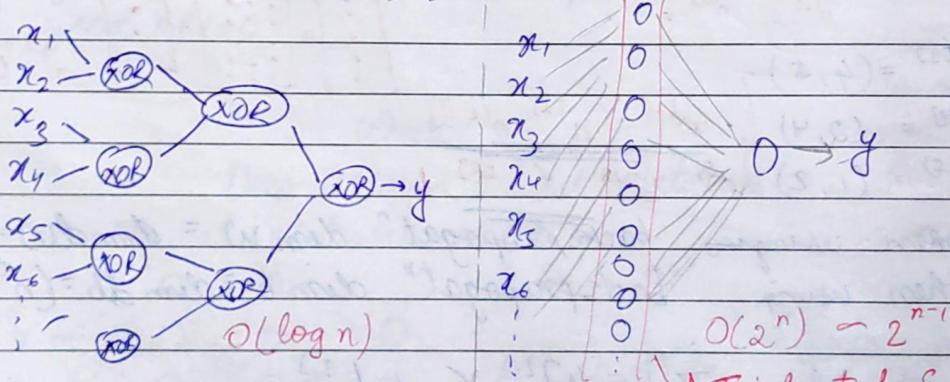
Ex for Audio detection :-

first layer might catch audibility of sound - high or low
second can be phonems - cat → b
after may be learn to recognize word & audio
finally recognise sentence & sounds / phrases

⇒ Circuit Theory & deep learning :-

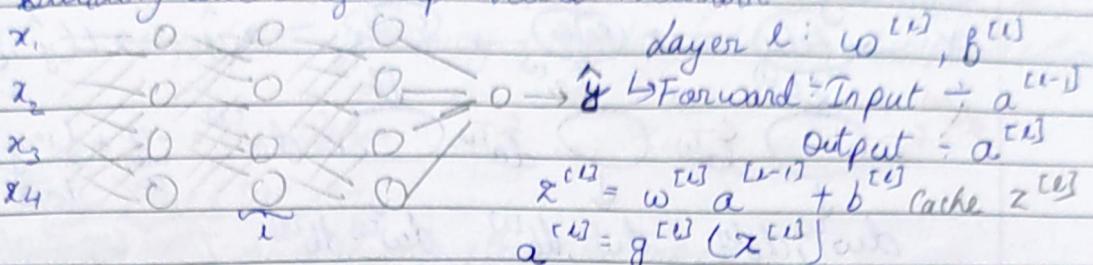
Informally:- There are functions you can compute with a small 1-layer deep neural network that shallower networks require exponentially more hidden units to compute.

Ex $x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } \dots \text{ XOR } x_n$

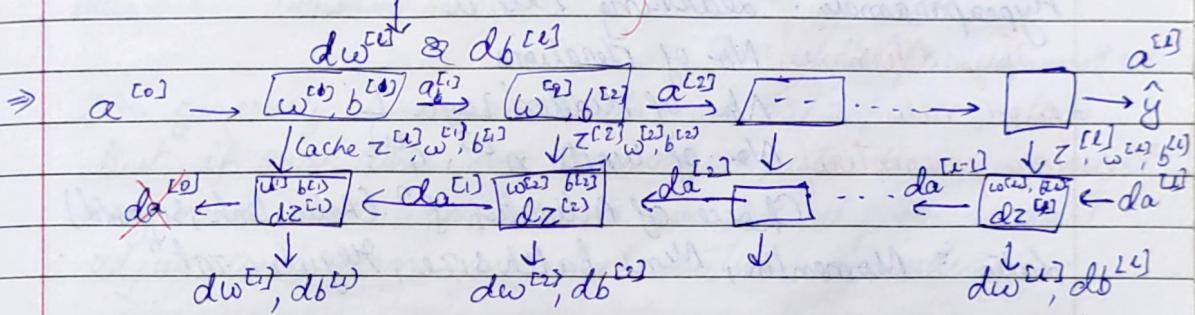
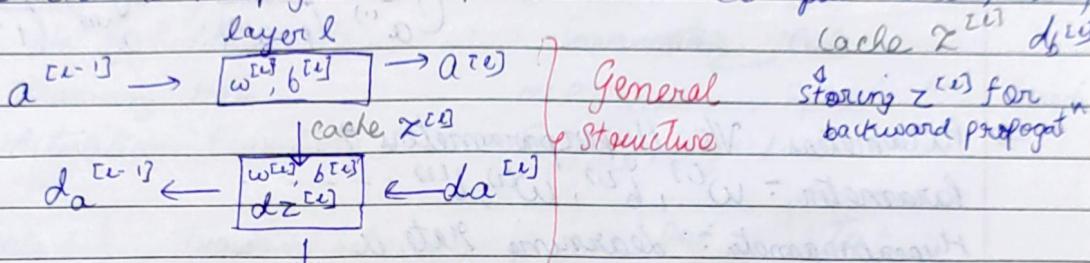


→ This has to be Exponentially large.

⇒ Building Blocks of Deep Neural Network :



↳ Backward Propagation Input = $da^{[l]}$ & Output = $da^{[l-1]}$, $dw^{[l]}$, $db^{[l]}$



$$w^{[l+1]} = w^{[l]} - \alpha dw^{[l]} \quad \& \quad b^{[l+1]} = b^{[l]} - \alpha db^{[l]}$$

⇒ Forward & Backward propagation :

i) Forward propagation for layer l

↳ Input $a^{[l-1]} \rightarrow [w^{[l]}, b^{[l]}]$

↳ Output $a^{[l]}$, Cache ($z^{[l]}$)

$$Z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \quad \left| \begin{array}{l} Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} = g^{[l]}(Z^{[l]}) \end{array} \right.$$

ii) Backward propagation for layer l :

↳ Input = $da^{[l]}$ & Output $da^{[l-1]}, dw^{[l]}, db^{[l]}$

$$dz^{[l]} = da^{[l]} \times g^{[l]}(Z^{[l]})$$

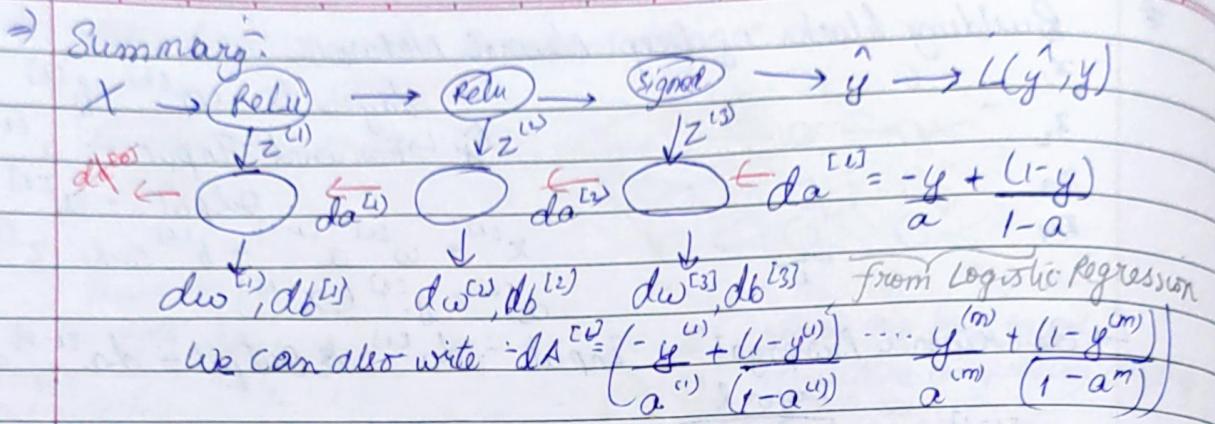
$$dw^{[l]} = dz^{[l]} \cdot a^{[l-1]T}$$

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

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

$$dz^{[l]} = W^{[l+1]T} \cdot dz^{[l+1]} \star g^{[l+1]}(Z^{[l]})$$

$$\left| \begin{array}{l} dZ^{[l]} = dA^{[l]} \star g^{[l]}(Z^{[l]}) \\ dw^{[l]} = \gamma m dZ^{[l]} \cdot A^{[l-1]T} \\ db^{[l]} = \frac{1}{m} \cdot np \cdot \text{sum}(dZ^{[l]}, \text{axis}=1, \text{keepdims=True}) \\ dA^{[l-1]} = W^{[l]T} \cdot dz^{[l]} \end{array} \right.$$



⇒ Parameters Vs Hyperparameters :-

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

Hyperparameter : Learning rate α

They determine No. of iterations

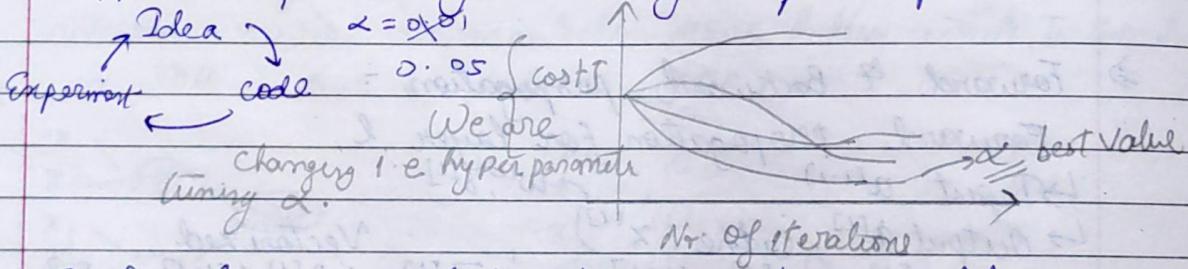
Final value No. of hidden layer L

using parameter No. of units $n^{[1]}, n^{[2]}, \dots$

Choice of Activation f^n (ReLU, tanh, Sigmoid)

Later : Momentum, Mini batch size, Regularization

⇒ Applied deep learning is a very empirical process :-



Ex Deep learning used in Vision, speech, NLP, Ad, Search, Recommendation.

DEEP LEARNING : Course - 2

Improving Deep Neural Networks :
Hyperparameter Tuning, Regularization and Optimization

Week - 1

- * Train / dev / test sets :

Applied ML is a highly iterative process -

- ↳ No. of layers
- ↳ No. of hidden units
- ↳ Learning rate
- ↳ Activation function

↑ Idea

Experiment ↘ Code

NLP, vision, speech, NLP, Search
Security, Logistics

Data :	Training Set	Dev	Test
--------	--------------	-----	------

Cross Validation / Development set

We generally divide our data in 80/20/20% ratios

But if our data has 1000,000 inputs then we might only need 10,000 inputs as our dev set. So u must adjust and plan your distribution accordingly.

- ⇒ Mismatched train/test distribution :

Training set :

Dev / test sets

Cat pictures from webpage

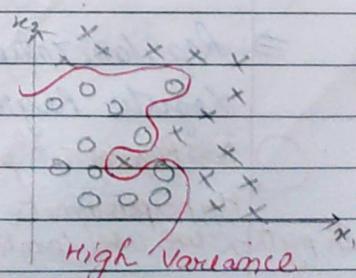
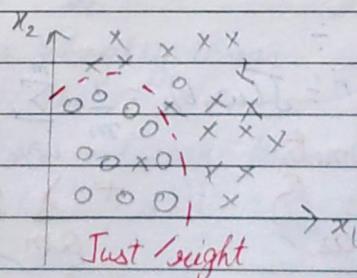
Cat pictures from users using your app

- ↳ Make sure dev & test come from same distribution.

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

Train / dev ^{test}

- ⇒ Bias / Variance :



Underfitting

Overfitting

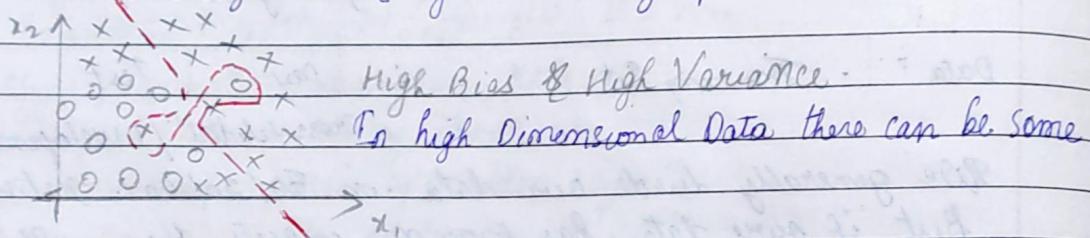
⇒ Bias and Variance:

Cat Classification	$y = 1/0$				
If train set error has	1%	15%	15%	0.5%	
Dev. set error has	11%	16%	30%	1%	

High variance, High bias & Variance, Low bias & low variance

Acc to humans we can easily classify it as cat or non-cat with 0% error. So on comparison 15% & 16% gives us high bias and vice versa.

⇒ How does high bias & high variance graph look like?



⇒ Basic "recipe" for machine learning:

High Bias?
(Training Data performance) → Bigger Network
Y (Train longer) [Diff NN architecture]

High variance?
(Dev Set performance) → More Data / regularization
Y (NN architecture search)

⇒ In early day we didn't had tools to improve any of bias or variance but with the help of deep learning we can individually improve bias or variance.

⇒ Regularization:

$$\text{Logistic Regression} : J(\omega, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) ; \omega \in \mathbb{R}^n, b \in \mathbb{R}$$

⇒ Regularizatⁿ parameter

$$+ \frac{\lambda}{2m} \|\omega\|_2^2 \quad \left. \begin{array}{l} \text{Regularizatⁿ term} \\ + \frac{\lambda}{2m} b^2 \end{array} \right\} \text{omit as } b \text{ is}$$

In python we use (lambda) as

$$\text{L2 Regularization} : \|\omega\|_2^2 = \sum_{j=1}^n \omega_j^2 = \omega^T \omega \quad \text{just one number}$$

L, regularization" $\frac{\lambda}{2m} \sum_{l=1}^{n_x} \|w\| = \frac{\lambda}{2m} \|w\|_F^2$

L, Not used much as w will be sparse and it will compress our data as it has a lot of zeroes.

\Rightarrow Neural Network:

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

$$\|w\|_F^2 = \sum_{i=1}^{n^{[L-1]}} \sum_{j=1}^{n^{[L]}} (w_{ij}^{[L]})^2 \text{ as } w: (n^{[L-1]}, n^{[L]})$$

\hookrightarrow Frobenius Norm: $\| \cdot \|_F^2 \rightarrow \| \cdot \|_F^2$

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

$$w^{[L]} = w^{[L]} - \alpha dW^{[L]} \quad \text{or}$$

$$\frac{\partial J}{\partial w^{[L]}} = dw^{[L]}$$

$\alpha \lambda$ regularization is also called as weight decay.

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

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

$$w^{[L]} := w^{[L]} \left(1 - \frac{\alpha \lambda}{m} \right) - \alpha (\text{from backprop})$$

We are subtracting $\frac{\alpha \lambda}{m}$ from $w^{[L]}$ so it got its name as weight decay.

\Rightarrow How does regularization prevent overfitting?

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

If $\lambda \uparrow$ then $w^{[l]} \downarrow$

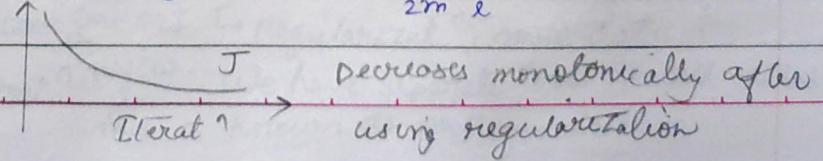
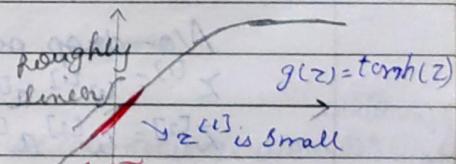
$$\downarrow \rightarrow x^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

$g(x)$ will be roughly linear

As we know, that if the f 's

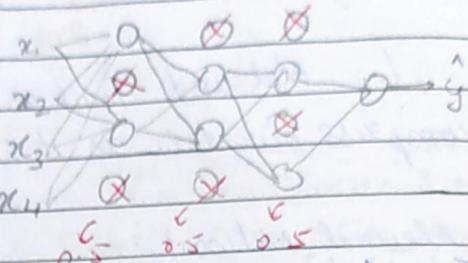
linear so the whole network becomes linear

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



⇒ Dropout Regularization :

$$\begin{array}{l} x_1 \quad 0 \quad 0 \quad 0 \\ x_2 \quad 0 \quad 0 \quad 0 \quad \rightarrow y \\ x_3 \quad 0 \quad 0 \quad 0 \\ x_4 \quad 0 \quad 0 \quad 0 \end{array}$$



we drop say (0.5) no. of nodes from a given layer and train a comparative small network.

⇒ Implementing dropout ("Inverted dropout")

Illustrate with layer $l=3$

$$d^3 = np.random.rand(a3.shape[0], a3.shape[1]) < \text{keep-prob}$$

Say keep-prob = 0.8 ; 0.2 chance of eliminating any hidden unit.

$$a3 = np.multiply(a3, d3) \quad i.e. a3 * = d3 [element wise]$$

Inverted dropout
ensure
 $a3 / = \text{keep-prob}$

expected value
Say we have 50 units $\therefore \sim 10$ units shut off

$$z^{[4]} = w^{[4]} a^{[3]} + b^{[4]} \quad i.e. \text{they are zero}$$

of a_3 remains

↳ Reduced by 20% ; 20% of $a^{[3]}$ = 0

the same

then we $/ = 0.8$ to not change value of $z^{[4]}$
it will bump back up roughly by 20% that we need

⇒ Making prediction at test time :

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

No drop out

$$x^{[0]} = w^{[1]} a^{[0]} + b^{[1]} ; a^{[1]} = g^{[0]}(z^{[0]})$$

$$x^{[2]} = w^{[2]} a^{[1]} + b^{[2]} ; a^{[2]} = \dots$$

\vdots \uparrow

⇒ Understanding dropout : why does drop-out work?

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

\rightarrow shrinking squared norm of the weights

\rightarrow Due to drop out we don't give any particular feature emphasis, we distribute it amongst all the other features.

$$\text{for } w^{[1]} = (3, 7); w^{[2]} = (7, 7)$$

$$w^{[1]} = (7, 3); w^{[2]} = (3, 2) \dots$$

x_1 0 0 0 0 0 \hat{y}
 x_2 0 0 0 0 0 \hat{y}
 x_3 0 0 0 0 0 \hat{y}

$\therefore w^{[2]}$ is biggest weight matrix
 → no overfitting and no dropout.

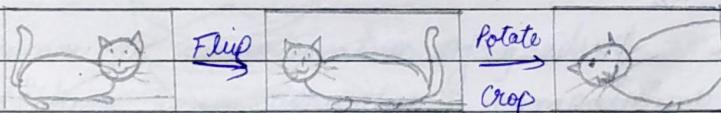
3. \rightarrow We will keep keep-prob relatively low say 0.5

In computer-vision we often use keep-prob.

Downside of dropout = $J(\cdot)$ is not defined in every iteration we are randomly calling out bunch of nodes.

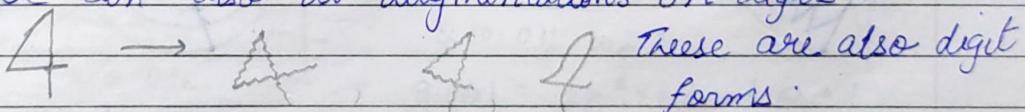
⇒ Other regularization methods :

Data augmentation



Using them we can add more data to our training set in an inexpensive way and reduce overfitting.

We can also do augmentations on digits:



⇒ Early Stopping : Due to early stopping we can not longer work simultaneously

① ⇒ Optimize cost $f(J)$ we have $\{J(w, b)\}$
 ↓
 midsize $l(w, b)$
 ↓
 dev set error \rightarrow G.D., Momentum, Adam, RMSprop

② ⇒ Not overfit -

train error or J \rightarrow Regularization, more data

$w \neq 0$
 No of iterat ↑ large w We have separate tools for both the task in NL
 this is known as orthogonalisation.

So instead of early stopping we can use L2 regularization but for that we have to try a lot of different values of λ , which can be computationally expensive.

\Rightarrow Normalizing inputs :

$$\begin{aligned} \text{Subtract Mean: } \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \\ \therefore x := x - \mu \end{aligned}$$

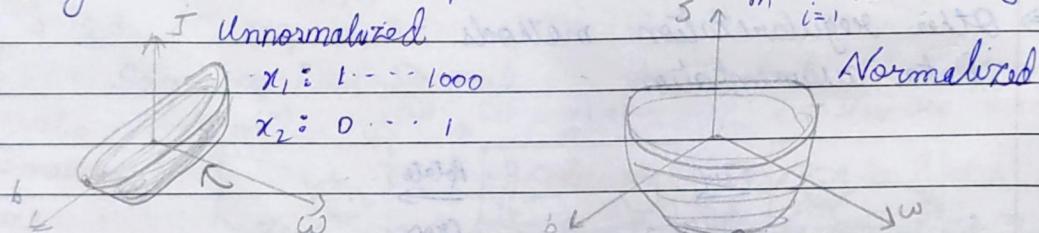
High variance
Low variance

$$\begin{aligned} \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2 \\ \therefore x' = \frac{x}{\sigma^2} \end{aligned}$$

Normalize variance

Use same σ^2 & μ to normalize test set.

\Rightarrow why to normalize inputs ? $J(w, b) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}^{(i)}, y^{(i)})$



\Rightarrow Vanishing / Exploding gradient :

$$x_1 \xrightarrow{w} y_1 \quad x_2 \xrightarrow{w} y_2 \quad \dots \quad \hat{y} = \underbrace{w^{[L]} \cdot w^{[L-1]} \cdots w^{[2]} w^{[1]}}_{\text{linear } f^n} \cdot x$$

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

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

$1.5^{L-1} x$

$$\begin{aligned} z^{[L]} &= w^{[L]} x \\ a^{[L]} &= g(z^{[L]}) = z^{[L]} \\ a^{[L]} &= g(z^{[2]}) = z^{[2]} = g(w^{[2]} a^{[1]}) \end{aligned}$$

If we take 0.5 instead of 1.5

we have 0.5 \times

for a $\sum w^{[l]} \geq 1$; "activation" can explode
deep network $w^{[l]} \leq 1$; activation will decrease exponential

⇒ Weight initialization for deep networks:

Single neuron example

$$x_i \rightarrow a^{[L]}$$

$$x_2 \rightarrow \hat{y}$$

$$x_3 \rightarrow a = g(z)$$

$$x_4 \rightarrow z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b^0$$

larger $n \rightarrow$ smaller w_i

$$\text{Variance } (w_i) = 1/n$$

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

↳ For ReLU activation $f(x)$

we see that:

We take variance as $2/n$

This is according to working

$$g^{[L]}(z) = \text{ReLU}(z)$$

↳ for tanh we use $1/n$

$$\text{Xavier initialization} \Rightarrow \sqrt{\frac{2}{n^{[L-1]} + n^{[L]}}}$$

⇒ Numerical approximation of gradient:

Checking your derivative computation

$$f(\theta) = \theta^3 \quad f'(\theta) = 3\theta^2$$

$$\frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon} \approx g(\theta)$$

$$(1.01)^3 - (0.99)^3 = 3.000$$

$$2(0.01)$$

$$g(\theta) = 3\theta^2 = 3 \text{ at } \theta = 1$$

$$\therefore \text{Apprx error} = 0.0001$$

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon} \rightarrow O(\epsilon^2) \quad ? \text{ More Accurate}$$

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} \rightarrow O(\epsilon) \quad ? \text{ Less Accurate}$$

⇒ Gradient checking : For a neural network

↳ Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ and reshape into a big vector θ and then concatenate

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

↳ Take $d\theta^{[1]}, d\theta^{[2]}, \dots, d\theta^{[L]}$ and reshape into a big vector $d\theta$ and then concatenate

Is $d\theta$ the gradient of $J(\theta)$

$$J(\theta) = J(\theta_1, \theta_2, \dots)$$

for each i :

$$d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\approx d\theta[i] = \frac{\partial J}{\partial \theta_i} \quad | \begin{array}{l} d\theta_{\text{approx}} \approx d\theta \\ \text{same dimension} \end{array}$$

Check $\|d\theta_{\text{approx}} - d\theta\|_2$; $\epsilon = 10^{-7}$

$\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2 \rightarrow$ If this value $\approx 10^{-7}$ or smaller
 great
 $10^{-5} \rightarrow \text{OK OK}; 10^{-3} \rightarrow \text{Wrong}$

⇒ Gradient checking for implementing notes :

↳ Don't use in training - only to debug [$d\theta_{\text{app}}[i], d\theta$]

↳ If algo fails grad check, look at components to try to identify by.

↳ Remember regularization $J(\theta) = \frac{1}{m} \sum L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum \|w^{(i)}\|_F^2$

$d\theta = \text{grad of } J \text{ w.r.t } \theta$

↳ Doesn't work with dropout [eliminate diff. random sets]
 or else keep prob=1.0 and then grad check

↳ Run at random initialization; perhaps again after
 some training : $w, b \leftarrow 0$

Week - 2

* Optimization Algorithms -

⇒ Mini - batch Gradient Descent -

Batch vs mini - batch gradient descent.

Vectorization allows you to efficiently compute on m examples

$$X = [x^{(1)} \ x^{(2)} \ \dots \ x^{(m)}] \rightarrow (n_x, m)$$

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

What if $m = 5,000,000$?

mini - batches of 1,000 each $[x^{(1)} \ x^{(1000)} \ x^{(2000)} \ \dots \ x^{(m)}]$

Same for y

Mini - batch t : x^{t+3}, y^{t+3}

x^{t+3}, y^{t+3} for mini batches

$$x^{t+3} \rightarrow (n_x, 1000) \quad y^{t+3} \rightarrow (1, 1000)$$

⇒ Mini - batch GD :-

for $t = 1, \dots, 5000$

Forward pass on x^{t+3}

This is also
called as

"epoch" one

pass through
the training set

$$\begin{aligned} z^{[t]} &= x^{[t]} X^{t+3} + b^{[t]} && \text{Vectorized implementation} \\ A^{[t]} &= g^{[t]}(z^{[t]}) && (\text{1000 examples}) \end{aligned}$$

$$A^{[t]} = g^{[t]}(z^{[t]}) \quad \rightarrow \text{for } x^{t+3}, y^{t+3}$$

$$\text{Compute cost } J = \frac{1}{1000} \sum_{i=1}^L \ell(\hat{y}^{(i)}, y^{(i)}) + \lambda \sum_{j=1}^{n_w} \|w^{[j]}\|^2$$

Backprop to compute gradients wrt J^{t+3} (using x^{t+3}, y^{t+3})

$$w^{[t]} := w^{[t]} - \alpha \nabla J^{t+3}; \quad b^{[t]} := b^{[t]} - \alpha \nabla J^{t+3}$$

⇒ Training with mini - batch GD :- There is noise as we are

Batch GD



Batch GD

Mini - batch GD



always training in new

batches x^{t+3}, y^{t+3} there can be some

x^{t+3}, y^{t+3}

noise while iterating

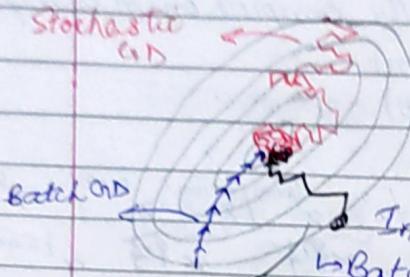
No. of iteratⁿ

mini batch iteratⁿ

⇒ Choosing your mini-batch size:

↳ If mini-batch size = m: Batch GP $(X^{fit}, Y^{fit}) = (X, Y)$

↳ If mini batch size = 1 Stochastic: Every Example is its own
 $(x^{(i)}, y^{(i)}) = (x^{(1)}, y^{(1)}) \dots$ mini batch



SGD won't even converge, always oscillate and wander around the region around the minimum.

In practice, mini-batch is in b/w 18 m

↳ Batch Gradient Descent - mini-batch size = m (Too long per iteration)

↳ SGD → mini-batch size = 1 Loose speedup from Vertoucat

In between mini-batch size (not too small not too big)

The gives us fastest learning
is we do get vectorized

"ii) makes progress without preceding entire train set

Choosing your ^{optimal} batch size

If small training set : Use batch GD ($m \leq 2000$)

Typical mini-batch size

64, 128, 256, 512 i.e powers of 2

Make sure mini-batch fits in CPU/GPU memory $X^{(k)}_{\text{fit}}$

\Rightarrow Exponentially weighted averages:

Temperature in London

$$\theta_1 = 40^\circ F \quad 4^\circ C \quad \text{Now}; \quad V_0 = 0$$

$$\theta_2 = 49^\circ F \quad 9^\circ C, V_i = 0.9V_0 + 0.1\theta_1$$

$$\theta_3 = 45^\circ F \quad ; \quad v_2 = 0.9 v_1 + 0.1 \theta_2$$

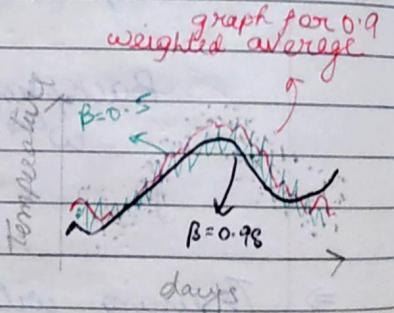
$$\theta_{180} = 60^\circ F \quad 15^\circ C$$

$$\theta_{B_1} = 56^{\circ}F \quad ; \quad v_f = 0.9v_{f-1} + 0.1\theta,$$

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

Kognitiv / Reaktion

moving average



$\rightarrow \beta = 0.9$ here ≈ 10 days temp

$\beta = 0.98$: $\sqrt{50}$ days temp

$$(1-\beta)\Theta_t \mid V_t = \perp ?$$

$$F_{cl} = 1 - \beta, \quad \beta = 0.98$$

$$= \frac{1}{e}$$

V_t as approximately average over $\frac{1}{1-\beta}$ days temperature.

For $\beta = 0.98$ we are getting a smoother curve as we taking average of last 50 terms
for $\beta = 0.5$: ≈ 2 days

→ Understanding Exponentially Weighted averages :

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

$$V_{100} = 0.9 V_{99} + 0.1 \theta_{100} \quad V_{100} = 0.9 V_{99} + 0.1 \theta_{100}$$

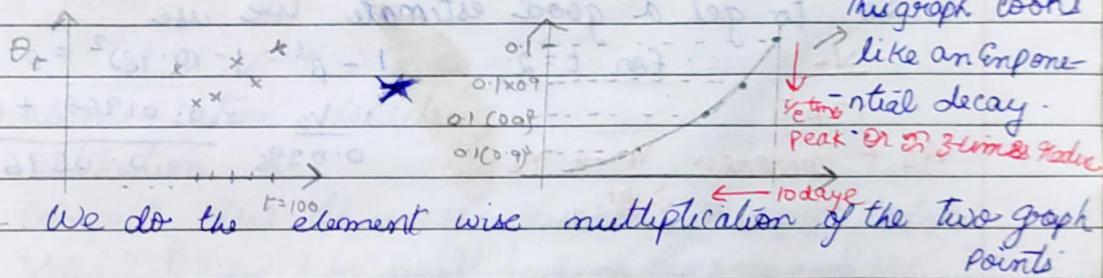
$$V_{99} = 0.9 V_{98} + 0.1 \theta_{99} \quad V_{100} = 0.1 \theta_{100} + 0.9(0.1 \theta_{99} + 0.9 V_{98})$$

$$V_{98} = 0.9 V_{97} + 0.1 \theta_{98}$$

....

$$\Rightarrow V_{100} = 0.1 \theta_{100} + 0.1 \times 0.9 \theta_{99} + 0.1 \times (0.9)^2 \theta_{98} + 0.1 (0.9)^3 \theta_{97} \dots$$

So what we are actually doing here



Some algebraic & calculus insights $0.9^{10} \approx 0.35 \approx \frac{1}{e}$

We know $(1-\epsilon)^{1/e} = \frac{1}{e}$; for $\epsilon = 0.1$

So we can say that around 10 days the peak θ_t will decay by $\frac{1}{e}$ times.

For different β we have diff parameters

Ex. $\beta = 0.98$; $0.98^{50} \approx \frac{1}{e}$; here it will take 50 days.
 $\epsilon = 1 - \beta$

⇒ Implementing exponentially weighted averages :

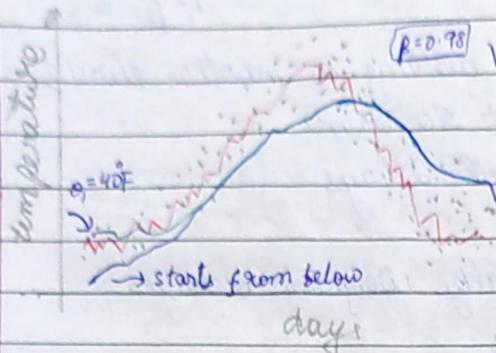
$$V=0$$

$V=0$ ← only variable we just substitute the other V_1, V_2, V_3, \dots so

$V_t = \beta V + (1-\beta) \theta_t$, Repeat! get next θ_t memory allocation is

$V := \beta V + (1-\beta) \theta_2$, $V := \beta V_0 + (1-\beta) \theta_1$ } very less.

⇒ Bias correction in exponentially weighted average:



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

$$V_0 = 0$$

$$V_1 = 0.98 V_0 + 0.02 \theta_1 \quad \text{---} \quad \text{D}$$

if $\theta_1 = 40^\circ F \therefore$ Using D

$$V_1 = 8^\circ F \text{ as } V_0 = 0 \text{ so}$$

we get a lower value.

On using "Bias correct" we get this blue curve and it starts from below

$$V_2 = 0.98 V_1 + 0.02 \theta_2$$

$$V_2 = 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2$$

$$V_2 = 0.0196 \theta_1 + 0.02 \theta_2 \quad [\because V_2 \text{ is also small?}]$$

It also isn't a very good estimate.

So to get a good estimate we use

$$\frac{V_t}{1-\beta^t} \quad \text{for } t=2; \quad 1-\beta^t = 1-(0.98)^2 = 0.0396$$

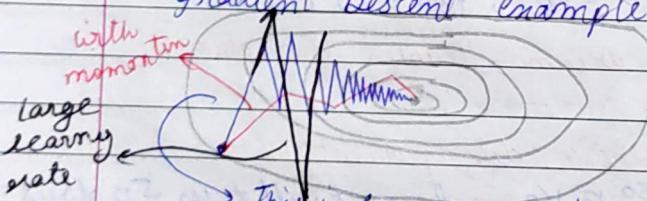
As t increases $V_t \rightarrow \bar{\gamma}_t$

$$\frac{V_t}{1-\beta^t}$$

$$\frac{V_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396} = \frac{0.0396}{0.0396}$$

⇒ Gradient Descent with momentum:

gradient descent example



This is how GD works
we must keep small learning rate

Or else, we might overbound.

What we want is that?

↓ slow learning ← fast learning

So we have

Momentum :

on iteration t:

$$V_{dw} = \beta V_{dw} + (1-\beta) d_w$$

$$V_{db} = \beta V_{db} + (1-\beta) d_b$$

$$w = w - \alpha V_{dw}$$

$$\beta = \beta - \alpha V_{db}$$

⇒ Adam Optimization Algorithm : It is basically RMSprop + Momentum
 $V_{dw} = 0$; $S_{dw} = 0$, $V_{db} = 0$, $S_{db} = 0$

On iteration t :

Compute dL/dw , dL/db using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dL/dw \quad \text{"momentum"} \quad \beta_1$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) dL/db$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) (dL/dw)^2 \quad \text{"RMSprop"} \quad \beta_2$$

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

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{1 - \beta_1 t}; \quad V_{db}^{\text{corrected}} = \frac{V_{db}}{1 - \beta_1 t}$$

$$S_{dw}^{\text{correct}} = \frac{S_{dw}}{1 - \beta_2 t}; \quad S_{db}^{\text{correct}} = \frac{S_{db}}{1 - \beta_2 t}$$

$$w := w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{correct}}} + \epsilon} \quad \text{&} \quad b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{correct}}} + \epsilon}$$

⇒ Hyperparameters choice :

α : needs to be tuned

β_1 : 0.9 (dL/dw)

β_2 : 0.999 ($(dL/dw)^2$)

ϵ : 10^{-8}

Adam = Adaptive moment estimation

⇒ Learning Rate decay :

1 epoch = 1 pass through the data

$$x^{1/3} | x^{2/3} | \dots$$

→ epoch 1
→ epoch 2

$$\alpha = \frac{1}{1 + \text{decay rate} * \text{epochnum}} \times \alpha_0$$

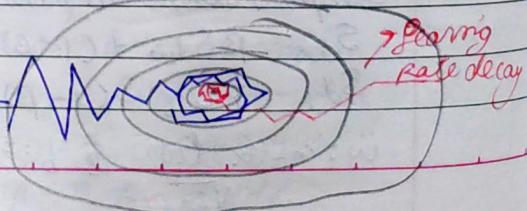
with this decay we can

easily approach our

minimum & get good

result

Big learning rate
converging is
not close to
minimum



Epoch	α	for $\alpha_0 = 0.2$
1	0.1	decay rate = 1
2	0.067	
3	0.05	α decreases
4	0.04	
⋮	⋮	

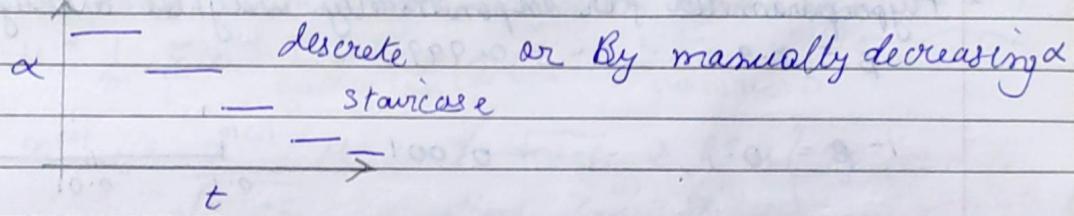
⇒ Other learning rate decay methods :-

different formulas

$$\alpha = 0.95^{\text{epoch num}} \quad \alpha_0 \leftarrow \text{Exponential decay}$$

\curvearrowright It can be any number less than 1

$$\alpha = \frac{k}{\sqrt{\text{epoch num}}} \quad \alpha_0 \text{ or } \frac{k}{\sqrt{t}} \alpha_0$$



Week - 3

* Hyperparameter tuning

⇒ Tuning process :-

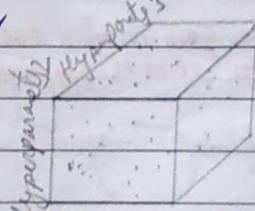
We have hyperparameters like $\alpha, B, R_1, R_2, \epsilon, \text{layers}, \text{hidden units}, \text{mini-batch size}, \text{learning rate decay}$

⇒ Try random values :- Don't use a grid

We can make a grid of numbers and select

$\alpha \in \dots$
Hyperparameters
Best pair or group

Hyperparameters!
more points.



○ → most important
○ → they are important too

For multiple hyperparameters we can use 3-D graph and select the best fit.

Here we now take smaller square and sample

⇒ Using an appropriate scale to pick Hyperparameters:

Picking Hyperparameters at random:

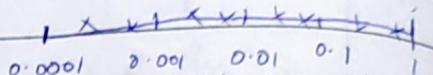
$$n^{(i)} = 80, \dots, 100$$

No. of layers L : 2 to 4 or 2, 3, 4

Appropriate scale for hyperparameters:

$$\alpha = 0.0001, \dots, 1$$

we can randomly search from



This is a logarithmic distribution

∴ Python code:

$$x = -4 * np.random.rand() \in [-4, 0]$$

$$\alpha = 10^x$$

$$\leftarrow \alpha \in [10^{-4}, 1]$$

$$\therefore a = \log_{10} 0.0001 = -4$$

$$b = \log_{10} 1 = 0 \quad \therefore x \in [a, b] \text{ i.e } [-4, 0]$$

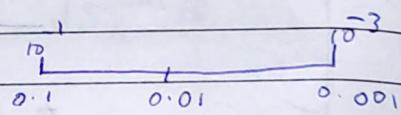
⇒ Hyperparameter for exponentially weighted averages:

$$\beta = 0.9, \dots, 0.999$$

Glauco 10 avg

↳ last 1000 avg

$$1 - \beta = 0.1, \dots, 0.001$$



$$\therefore \beta = 1 - 10^x$$

$$x \in [-3, -1]$$

• We must use a log scale not a linear scale.

$$1 - \beta = 10^x$$

⇒ Hyperparameter tuning in practice: Pandas Vs Caviar-

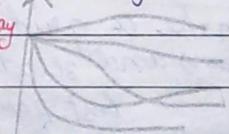
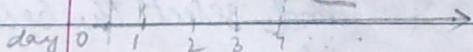
We keep on changing hyperparameters as intuitions do

get stale and re-evaluate your model occasionally.

Babysitting one model

Training many models in parallel

we keep a check on our model each day
and淘汰 the worst ones to select the hyperparameters



Caviar - Best model is selected
like best fishes are selected as
there are millions of eggs.

Panda - only have 1-2 babies
and it takes care of it

⇒ Batch normalization :-

Normalizing activations in a network :-

Normalizing inputs to speed up learning

$$x_1 \quad w, b \quad \text{---} \quad O \rightarrow \hat{y} \quad \mu = \frac{1}{m} \sum_i x^{(i)}$$

$$x_2 \quad \text{---} \quad \sigma^2 = \frac{1}{m} \sum_i x^{(i)2}$$

$$x_3 \quad \text{---} \quad X := X - \mu \quad \sigma^2 \rightarrow \text{elementwise}$$

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

Can we normalize $a^{[2]}$ so as to train $w^{[3]}, b^{[3]}$ faster

Normalize $x^{[2]}$

⇒ Implementing Batch Norm :-

Given some intermediate values in Neural Net $(x^{(1)} \dots z^{(m)})$

$$\mu = \frac{1}{m} \sum_i x^{(i)} ; \sigma^2 = \frac{1}{m} \sum_i (x^{(i)} - \mu)^2$$

$$x_{\text{norm}}^{(i)} = \frac{x^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \gamma = \sqrt{\sigma^2 + \epsilon}, \quad \beta = \mu \quad \left. \begin{array}{l} \text{If } \gamma \neq 0 \\ \text{then } \tilde{x}^{(i)} = x^{(i)} \end{array} \right\} \text{For diff components}$$

$$\tilde{x}^{(i)} = \gamma x_{\text{norm}}^{(i)} + \beta \quad \text{so that we don't always get variance = 1}$$

↳ learnable parameters of model & mean = 0

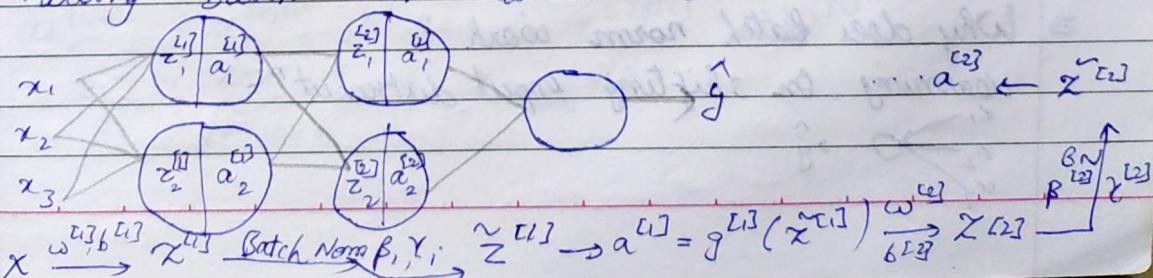
So we can $\tilde{x}^{(i)}$ instead of $x^{(i)}$

Because we use $\sigma^2 = 1 \& \mu = 0$ for our data

so our data will be on a linear fit

so we must choose different values of γ & β to adjust our model.

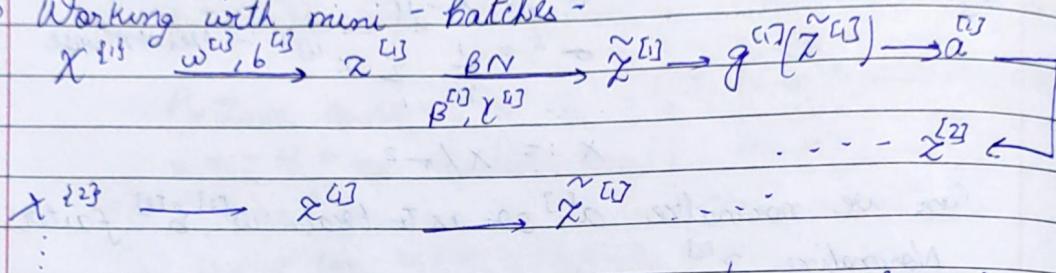
⇒ Fitting batch norm into a neural network :-



Parameters = $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$
 $\beta^{[1]}, \gamma^{[1]}, \dots, \beta^{[L]}, \gamma^{[L]}$ $\beta^{[L]} = \beta - \alpha d\beta^{[L]}$

This β & Adam's β is different.
 We can use "batch normalization" with one line of code
 i.e. `tf.nn.batch_normalization`.

⇒ Working with mini-batches:-



Parameters = $w^{[1]}, b^{[1]}, \beta^{[1]}, \gamma^{[1]}, \dots, w^{[L]}, b^{[L]}$ $X^{[t]} = w^{[t]} a^{[t-1]} + b^{[t]}$
 $\beta^{[t]}$ gets eliminated coz we do normalization on $X^{[t]}$ take its mean so $b^{[t]}$ doesn't have significant role. So we take $b^{[t]} = 0$ $X_{\text{norm}}^{[t]} = \tilde{X}^{[t]} = \gamma^{[t]} X_{\text{norm}}^{[t]} + \beta^{[t]}$

⇒ Implementing gradient descent:-

for $t=1 \dots \text{num_mini_batch}$

Compute forward prop on $X^{[t]}$

In each hidden layer use Batch Norm

to replace $X^{[t]}$ with $\tilde{X}^{[t]}$

Use backprop to compute $dw^{[t]}, db^{[t]}, d\beta^{[t]}, d\gamma^{[t]}$

Update parameters: $w^{[t]} := w^{[t]} - \alpha dw^{[t]}$

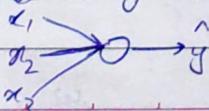
$\beta^{[t]} := \beta^{[t]} - \alpha d\beta^{[t]}$

$\gamma^{[t]} := \gamma^{[t]} - \alpha d\gamma^{[t]}$

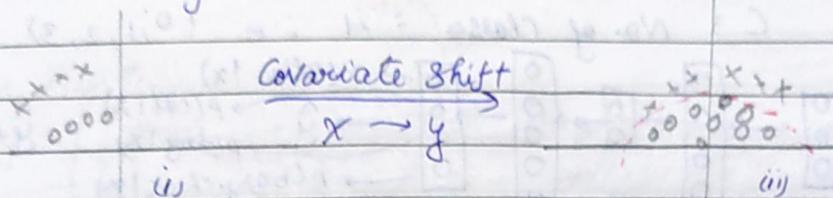
works with momentum; RMSprop, Adam.

⇒ Why does batch norm work?

Learning on shifting input distribution:-



When we classify and train our model to identify cat and non-cat; what if we only provide computer with Black cat images and then ask for classification. Our model will fail. In Fig (i) it is hard to draw a decision boundary.



So what does batch norm do; when operating with different batches we might get different results but mean and variance will always be the same on plotting it on graph.

\Rightarrow Batch norm as regularization:

- \hookrightarrow Each mini-batch is scaled by the mean / variance computed on just that mini-batch.
- \hookrightarrow This adds some noise to the values $x^{(i)}$ within that minibatch. So similar to dropout; it adds some noise to each hidden layer's activations.
- \hookrightarrow This has a slight regularization effect.

\Rightarrow Batch norm at test time:

$$\mu = \frac{1}{m} \sum_i x^{(i)} ; \sigma^2 = \frac{1}{m} \sum_i (x^{(i)} - \mu)^2 \quad m \text{ is no. of samples in mini-batch}$$

$$x_{\text{norm}}^{(i)} = \frac{x^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} ; \tilde{x}^{(i)} = \gamma x_{\text{norm}}^{(i)} + \beta$$

μ, σ^2 : estimate using exponential weighted average (across mini-batches)

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

$$\downarrow \quad \downarrow \quad \dots$$

$$\mu^{(1)} \quad \mu^{(2)} \quad \text{ & Same for } \sigma^2 \quad \sigma^2 \quad \dots$$

$$x_{\text{norm}} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} ; \tilde{x} = \gamma x_{\text{norm}} + \beta$$

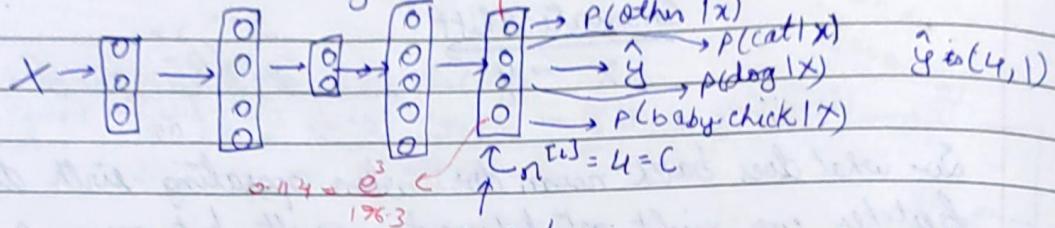
Multi Class Classification:

Softmax Regression:

Recognizing cats, dogs and baby chicks, other

$$\begin{matrix} 1 & 2 & 0.84 & e^5 \\ & & \frac{176.3}{e^5} & 3 \\ & & 0.114 & 0 \end{matrix}$$

$$C = \text{No. of Classes} = 4 \text{ i.e. } (0, 1, 2, 3)$$



$$x^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}$$

Activation function: Softmax f^n

$$t = e^{(z^{[L]})} \quad (4, 1)$$

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum t_i}$$

$$\frac{0.84}{e^5 + e^2 + e^{-1} + e^3}$$

$$0.114$$

(4, 1)

$$a^{[L]} = g^{[L]}(x^{[L]})$$

(4, 1)

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

$$\therefore \sum t_i = 176.3$$

$$\therefore \sum t_i = 176.3$$

$$a^{[L]} = \frac{t}{176.3}$$

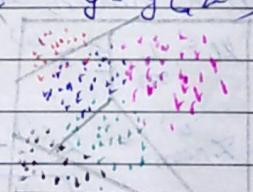
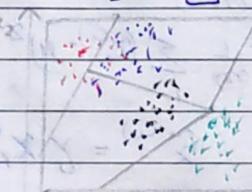
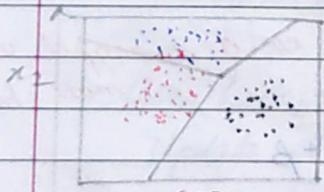
\Rightarrow Softmax Example:

$$x_1 \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \hat{y}$$

$$x_2 \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \hat{y}$$

$$x^{[L]} = w^{[L]} x + b^{[L]}$$

$$a^{[L]} = \hat{y} = g(x^{[L]})$$



\Rightarrow Training a softmax classifier: $C=4$

$$x^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}; t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}; g^{[L]}(x^{[L]}) = \begin{bmatrix} e^5 + e^2 + e^{-1} + e^3 \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.812 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

$g^{[L]}(\cdot) \rightarrow$ Softmax Activation f."

hard max

→ Softmax regression generalizes logistic Regression to C classes

If $C = 2$; softmax reduces to logistic Regression.
 $a^{(1)} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$ Single output

⇒ Loss Function:

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \text{ cat } \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} = a^{(4)} \text{ So our model is not performing well.}$$

$$\left[\begin{array}{l} y_1 = y_3 = 0 \\ y_2 = 1 \end{array} \right] \quad \alpha(\hat{y}, y) = -\sum_{j=1}^4 y_j \log \hat{y}_j \Rightarrow -y_2 \log \hat{y}_2 = -\log \hat{y}_2$$

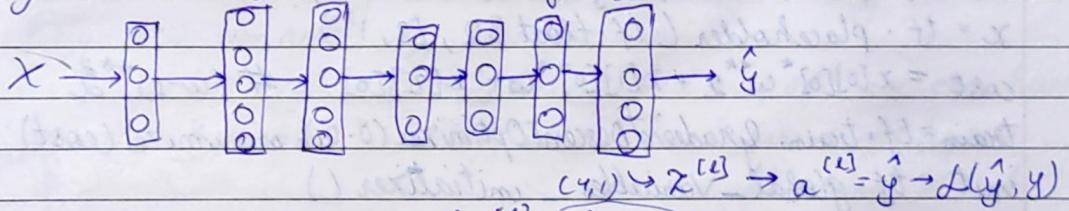
To reduce cost, make \hat{y}_2 as big as possible.

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

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

$$= \begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \rightarrow (4, m) \quad = \begin{bmatrix} 0.3 & 0.2 & 0.1 & 0.4 \\ 0.2 & 0.1 & 0.3 & 0.4 \\ 0.1 & 0.3 & 0.2 & 0.4 \\ 0.4 & 0.4 & 0.4 & 0.4 \end{bmatrix} \rightarrow (4, m)$$

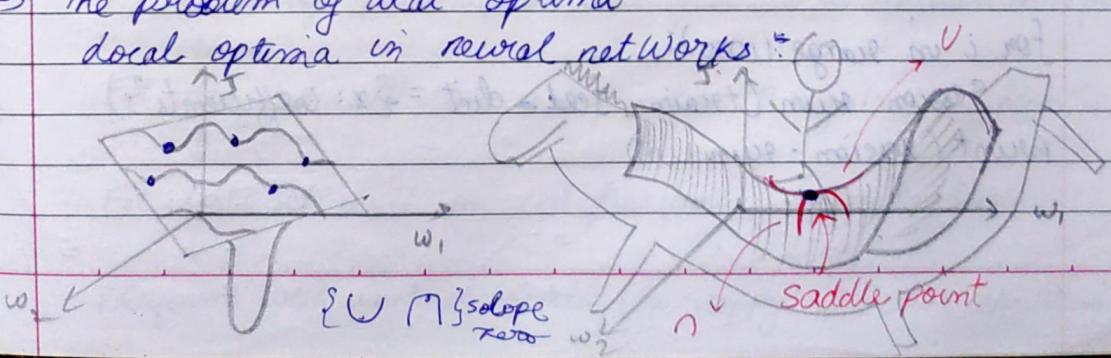
⇒ Gradient Descent with softmax:



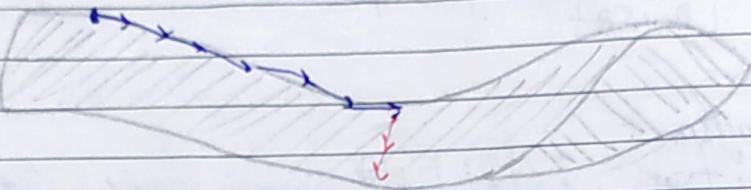
$$\text{Back propagation: } d\hat{y}^{(i)} = \hat{y} - y^{(i)} \quad \frac{\partial J}{\partial z^{(i)}}$$

⇒ The problem of local optima:

local optima in neural networks:



⇒ Problem of plateau :-



- ↳ unlikely to get stuck in a bad local optima
- ↳ Plateaus can make learning slow.

⇒ Tensorflow :-

Motivating problem

$$J(w) = w^2 - 10w + 25 = (w-5)^2 \text{ at } w=5 \text{ minimum}$$

Code Example :-

```
import numpy as np
import tensorflow as tf
coefficients = np.array([1, -20, 25])
w = tf.Variable([0], dtype=tf.float32)
x = tf.placeholder(tf.float32, [3, 1])
cost = x[0][0] * w**2 + x[1][0] * w + x[2][0] # (w-5)**2
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
init = tf.global_variables_initializer()
session = tf.Session()
```

$x[0][0] \rightarrow w^2 \Rightarrow [x[0][0] \cdot w^4]$ ← Cost →

Also with $tf.Session()$ as session:
written $Session.run(init)$
as $print(session.run(w))$

for i in range(1000):

```
    session.run(train, feed_dict={x: coefficients})
    print(session.run(w))
```

DEEP LEARNING : Course-3

Introduction to ML Strategy

Week - 1

★ Structuring ML Projects

*) Ideas :

↳ Collect more data

↳ Collect more diverse training set

↳ Train algo longer with Gradient Descent

↳ Try Adam instead of Gradient Descent

↳ Try bigger network

↳ Try smaller network

↳ Try dropout

↳ Add L2 regularization

↳ Network Architecture

- Activation functions

- No of hidden units

⇒ Orthogonalization :

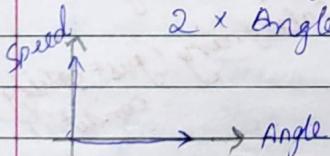
Car Tuning Example :

We have Steering, Acceleration, Breaking what if they are connected by this relation

$0.3 \times \text{Angle} - 0.8 \text{ Speed}$ } But instead we can have

$2 \times \text{Angle} + 0.9 \text{ Speed}$ } 2 different features

i.e Angle and Speed



⇒ Chain of Assumption in ML :

Fit training set well on cost function. (⇒ Human level perform)



Fit dev set well on cost function ⇒ Regularization, Bigger Training Set



Fit test set well on cost function ⇒ Bigger Dev Set



Performs well in real world (Ex: Happy cat picture app users) ⇒ cost function

⇒ Single number evaluation metric : of the examples recognised as cat, what are actually cats?

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

of all the images really are cats, what % were correctly recognised by your classifier?

$$F\text{-score} = \frac{2PR}{P+R} \quad \text{Harmonic mean of P \& R}$$

Dev set + Single real number evalutⁿ metric → speeds up & Improve this iterative process.

⇒ Another example : What if we are training a lot of models on different parameters (example : of train error).

Algorithm	US	China	India	Average	we can check
A	3%	7%	5%	5%	the performance
B	5%	6%	4%	5%	by taking average.
C	2%	3%	4%	3%	
D	5%	8%	2%	5%	

⇒ Satisfying and optimizing metrics :

Another cat classification example : → Optimizing (Max)

Classifier	Accuracy	Running Time	Satisfying
A	90%	80ms	Just satisfy
B	92%	95ms $\leq 100ms$	condition a
C	95%	1500ms	

Say cost = accuracy - $0.5 \times$ running time.

So we maximize accuracy & subject to running time $\leq 100ms$.

N metrics : 1 Optimizing

N-1 Satisfying

⇒ Wake words / Trigger words :

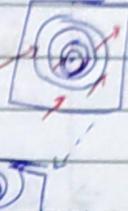
Alexa, OK Google, Hey Siri, nihubaidu (1% & 3%)

Accuracy & false positive → Maximize Accuracy such that ≤ 1 false positive every 24 hours.

3 Train / dev / test distributions :

Cat classification dev / test sets

US { Dev You must randomly shuffle all the data
 UK } wrong and then divide it into dev / test sets
 Europe method → If we chose this idea then
 India } Test we are particularly aiming at this
 China } target to gain accuracy but what
 Aus if target is shifted from the original
 } space our model will fail to obtain
 High level accuracy.


⇒ Guideline : same distribution

choose a dev set and test set to reflect data you expect to get in the future and consider important to do well on

⇒ Size of dev / test sets :

Old way of splitting data :

70%	30%	Good for 100 1000 10000 Examples
Train	Test	
60%	20% 20%	for 1,000,000 Examples
Train	Dev Test	
98%	1% 1%	
Train	Dev Test	

→ Set your test set to big enough to give high confidence on the overall performance of your system.

⇒ When to change Dev / Test sets and metrics :

Cat dataset Example :

Metric : Classification Error

✗ Algo A : 3% Error → It shows some pornographic images

✓ Algo B : 5% Error → Not showing any pornographic image

Metric + Dev → prefers A. // But You / user → prefers B.

$$\text{Error} \Rightarrow \sum_{\text{Error}} \sum_{i=1}^m P_i y_{\text{pred}}^{(i)} + y^{(i)} f_i w^{(i)}$$

↳ predicted value (0/1)

$$w^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is non-porn} \\ 100 & \text{if } x^{(i)} \text{ is porn} \end{cases}$$

So now our algo on iterating through pornographic images it will have a huge error.

→ Orthogonalization for cat pictures : Anti-porn

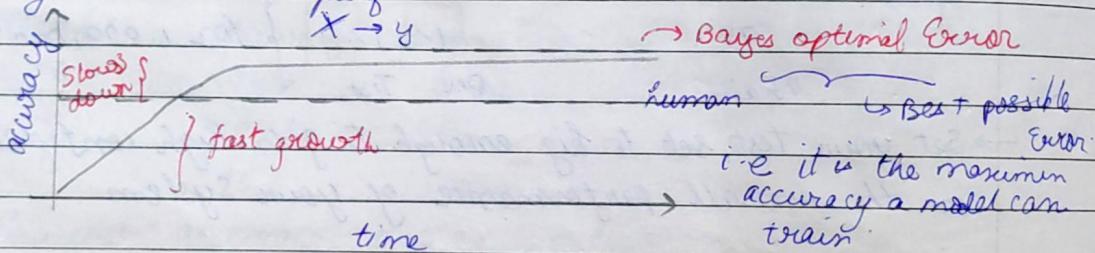
↳ So far we've only discussed how to define a metric to evaluate classifiers. ≈ place the target

↳ Worry separately about how to do well on this metric.
≈ It is a separate prob how to shoot the target.

$$J = \frac{1}{\sum w^{(i)}} \sum_{i=1}^m \omega^{(i)} \delta(\hat{y}^{(i)}, y^{(i)})$$

↳ If your algo/model is doing well on your metric + dev/test set does not correspond to doing well on your application, change your metric and/or dev/test set.
As user images can be blurry, rotated.

→ Why human-level performance :



Or cat images which are very blurry. Cor this can be some example say (audio) which is very noisy and can't be transcribed.

↳ As humans are also good at classifying cat images or audio so the gap is very short.

- Why compare to human-level performance?
 - Humans are quite good at a lot of tasks so as long as ML is worse than humans you can -
 - Get labeled data from humans (X, y)
 - Gain insight from an manual error analysis:
 - Why did a person get this right?
 - Better analysis of bias/variance:

⇒ Avoidable Bias:

Cat classification =

Bayes vs Humans	1%	7.5%	\uparrow	\leftarrow	Avoidable Bias; we can't go below
Training Error	8%	8.5%	\downarrow	\leftarrow	Bayes error
Dev. Error	10%	10.5%	\downarrow	\leftarrow	Variance
			Focus on ↓	Focus on ↓	
			Bias	Variance	

Human level error as a proxy for Bayes Error

⇒ Understanding human-level performance:

Human level error as a proxy for Bayes Error

Medical image classification Example

Suppose -

- | | | |
|--|------|--|
| a) Typical human - - - 3% error | 3% | What is human-level error?
Bayes error $\leq 0.5\%$. |
| b) Typical doctor - - - 1% error | 1% | |
| c) Experienced doctor - - - 0.7% error | 0.7% | |
| d) Team of exp. doc - - - 0.5% error | 0.5% | |

⇒ Error analysis Example:

Human (proxy for Bayes Error)	5%	13% \uparrow 0.5% ↓ 4.5% ↓ 0.1%	0.5%
↑ Avoidable bias			↓ 0.2% Available Bias
Training Error	- - - 5%	1%	0.7%
↓ Variance		↓ 4%	↓ 0.1%
Dev. Error	- - - 6%	5%	0.8%
Train a bigger network		Bias	Variance: Regularization?

⇒ Surpassing human level performance :-

Team of humans	0.5%	Avoidable Bias = 0.1% / 0.5%.	Is Avoidable
One human	1%	1%.	Bias = 0.1%.
Training Error	0.6%	Variance 0.2%.	0.2%.
Dev Error	0.8%	0.3%. As 0.3% might 0.4%. Be an overfit so incomplete data	0.3%.

⇒ Problems where ML significantly surpasses human level performance :-

↳ Online Advertising [Click on Ad] } lots of data

↳ Product Recommendations

↳ Logistics (predicting transit time)

↳ Loan approvals

↳ These are structured data ; humans are good at Natural perception problems like computer vision , audio , video not at structured data.

↳ Some image recognition } Computers have surpassed humans

↳ Some speech recognition }

↳ Some medical . Ex ECG

⇒ Improving your model performance :-

The two fundamental assumptions of supervised learning

↳ You can fit the training set pretty well. \Rightarrow Avoidable Bias

↳ The training set performance generalizes pretty well to the dev/test set. \Rightarrow Variance

⇒ Reducing (avoidable) bias and variance :-

Human-level \Rightarrow Train bigger model

Avoidable \rightarrow Train longer / better optimization algorithm Adam , RMSProp

Bias \rightarrow NN architecture / hyperparameter search (ANN, CNN)

Training Error { More data

Variance \rightarrow Regularization (L_2 , dropout, data augmentation)

Dev Error \rightarrow NN architectural hyperparameter search

Week - 2

* Error Analysis

Carrying out error analysis:

→ Should you try to make your cat classifier do better on dogs?

Error analysis

↳ Get ~100 mislabeled dev set examples.

↳ Count up how many dogs. "Ceiling Method"

If our model is showing 10% of error and if we have 5/100 dogs and we manually classify them. But then too our error will drop down to 9.5% from 10% error.

But if there are 50/100 dogs then our error will drop down to 5% from 10%.

→ Evaluate multiple ideas in parallel:

↳ Fix pictures of dogs being recognized as cats.

↳ Improve performance on blurry images.

↳ Fix great cats (Lions, panthers, etc.) being misrecognized.

Incorrectly labeled	Image	Dog	great cats	Blurry	Instagram comments	
✓	1	✓			✓	Pitbull
✓	2			✓	✓	
✓	3		✓	✓		Bailey day of 200
6/1	✓. total	8%	43%	61%	12%	

→ Cleaning up ~~incorrectly~~ labeled data:

x → Cat dog Cat cat Mat dog dog

y → 1 0 1 1 0 1 0

DL Algorithms are quite robust to random errors in the training set.

For classifying these incorrect labels we can add another column of incorrectly labeled data.

Now the question arises that whether we should fix these 6% labels or not. Before that we look at these 3 cases :-

- i) Overall dev/test error - - - 121.10%. So instead of solving Errors due to incorrect labels - - - 0.6%, 0.6%, 0.6% error we can Errors due to other causes - - - 104%, 9.4% focus on other problems
- ↳ Fix up dev/test errors here

→ Correcting incorrect dev/test set examples :-

- ↳ Apply same process to your dev and test sets to make sure they continue to come from the same distribution.
- ↳ Consider examining examples your algorithm got right as well as ones it got wrong.
- ↳ Train and dev/test data may now come from slightly different distributions.

→ Build your first system quickly then iterate :-

- ↳ Noisy Background
 - Café Noise, Car Noise
 - Set up dev/test set and metric
- ↳ Accented Speech
- ↳ Far from microphone
- ↳ Young children's speech
- ↳ Stuttering (uh, ah, hmm...)
- Build initial system quickly
- Use bias/variance analysis & error analysis to prioritize next steps.

→ Training and Testing on Different Distributions :-
Cat App example

Data from webpage
Sharp, Clear images

Data from mobile app
Blurry Images

~ 200,000 images → ↪ ~ 10,000 images

210,000 Shuffle

~~1st~~

option :-

Train

Dev

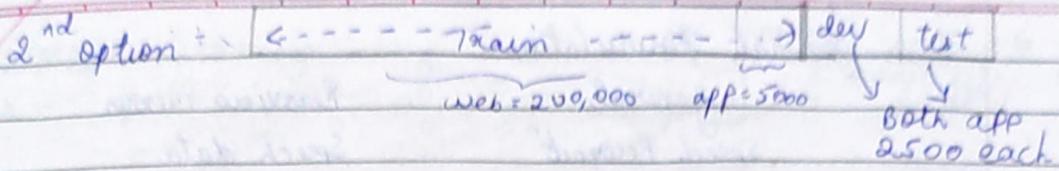
Test

→ 2500

205,000

→ 2500

→ 2381 - webpage
119 - mobile app



⇒ Speech Recognition example:

Ex: Speech activated rearview mirror

Training

Purchased data

Smart speaker control

Voice keyboard

----- ↗ 500,000

Dev/Test

Speech activated rearview mirror

↳ 20,000

Train

Dev Test

510K

5K 5K

10K

⇒ Cat classifier example:

Assume humans get 60% error, Train-dev error - 9% Problem

Training Error - 17% ↘ 19%

Dev error - 10%

Training Error - 1% ↗ Variance

Dev Error - 10% ↗ Problem

Dev Error - 10%

Train Dev Dev Dev

NN saw this

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

Training Error - 1% ↗ Bias ↗ Variance ↗ Available Bias

Training Error - 1% ↗ Bias ↗ Variance ↗ Available Bias

Train - dev Error - 1.5% ↗ data mismatch ↗ Bias ↗ Variance

Dev error - 10% ↗ data mismatch ↗ Bias ↗ Data mismatch

Bias ↗ Bias + data-mismatch

⇒ Bias Variance on mismatched training and dev/test sets:

Human level - 4% ↗ available bias ↗ Variance ↗ Bias ↗ Variance ↗ Training set evaluated

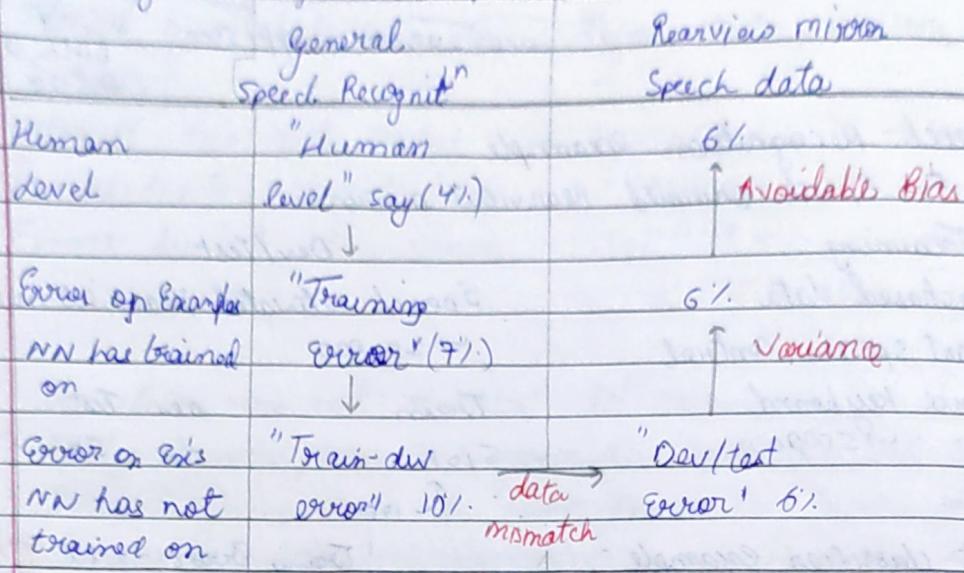
Training set error - 7% ↗ Variance ↗ Training set evaluated

Train - dev set error - 10% ↗ Data mismatch ↗ Dev/best set evaluated

Dev - error - 12% ↗ degree of overfitting to dev set ↗ Dev/best set evaluated

Test - error - 12% ↗ overfitting to dev set ↗ Dev/best set evaluated

⇒ More general formulation :-



⇒ Addressing Data Mismatch :-

→ Carry out manual error analysis to try to understand difference between training and dev/test sets.

Ex: noisy - car noise ; mismatching street numbers

→ Make training data more similar , or collect more data similar to dev/test sets.

Ex: Simulate noisy in-car data

⇒ Artificial data synthesis :-

It has all the alphabets $\boxed{1}) + \boxed{2}) = \boxed{3})$

$\boxed{1})$ "The quick brown fox jumps over the lazy dog"

↪ 10,000 hours

Car

noise

↪ 1 hour of
car noise

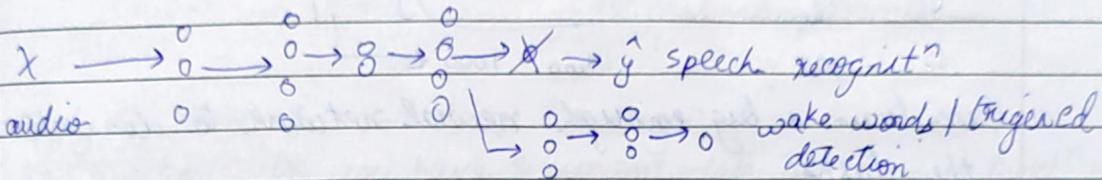
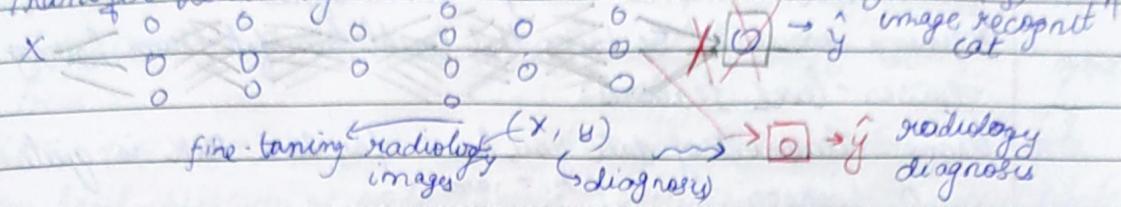
Synthesized

in-car audio

↪ Set of all audio
in car.

There is a chance that 1 hour of car noise can overfit the data so one ^{can} get 10,000 hours of car noise data.

⇒ Transfer Learning :



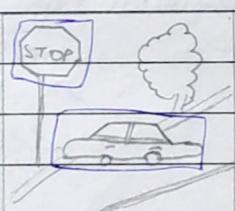
↳ When does transfer learning make sense?

Transfers from $A \rightarrow B$

- i) Task A and B have the same input X [like image, audio]
- ii) You have a lot more data for Task A than task B.
- iii) Low level features from A could be helpful for learning B.

⇒ Multi-task learning :

Simplified autonomous driving example :

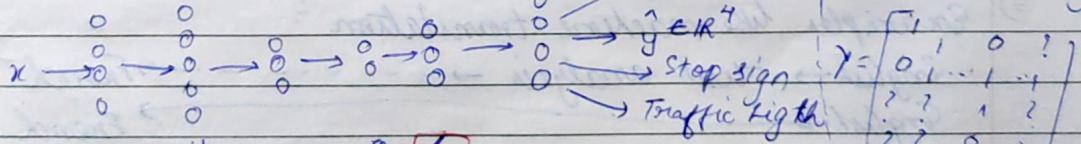


$x^{(i)}$	Pedestrian	0
	Cars	1
	Stopsigns	1
	Traffic light	0

$$y^{(i)} / (4, m)$$

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

⇒ Neural Network Architecture :



Loss $\hat{y}^{(i)}$

sum only over values of j with $0/1$ label

$$\frac{1}{m} \sum_{l=1}^m \sum_{j=1}^n \ell(\hat{y}_j^{(i)}, y_j^{(i)})$$

↳ usual length

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

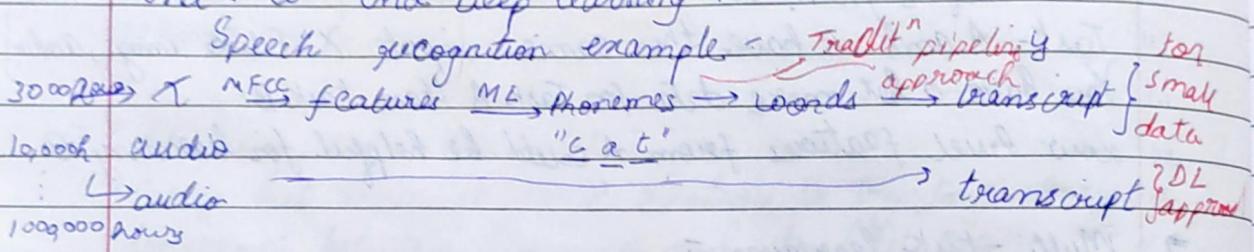
Unlike softmax regression - here one image can have multiple labels like cars, pedestrian ...

- ↳ Whether multi-task learning make sense?
- ↳ Training on a set of tasks that could benefit from having shared lower-level features.

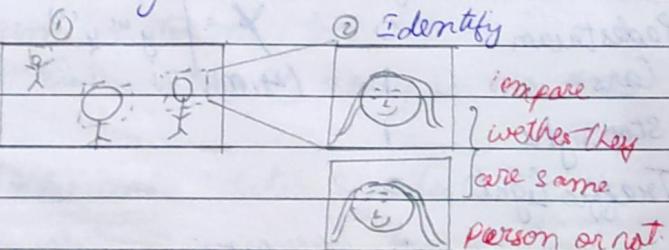
- ↳ Using : amount of data you have for each task is quite similar
- | | | | |
|---|-----------|------------------|--------------|
| A | 1,000,000 | $A_1 = 1000$ | $\} 99,000$ |
| B | 1000 | $A_2 = 1000$ | |
| | | $A_{100} = 1000$ | \leftarrow |

- ↳ Can train a big enough neural network to do well on all the tasks.

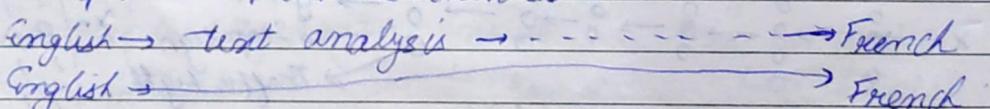
⇒ End-to-End Deep learning :



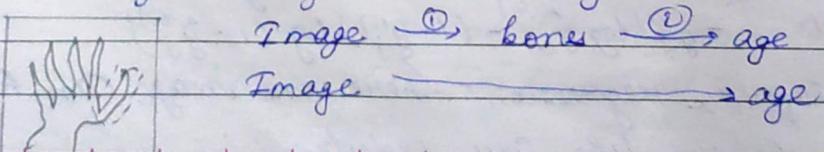
⇒ Face Recognition :



⇒ Examples like machine translation :



⇒ Estimating child age from X-ray :



⇒ Whether to use end-to-end learning:

Pros and cons of end-to-end deep learning

Pros =

Cons

input

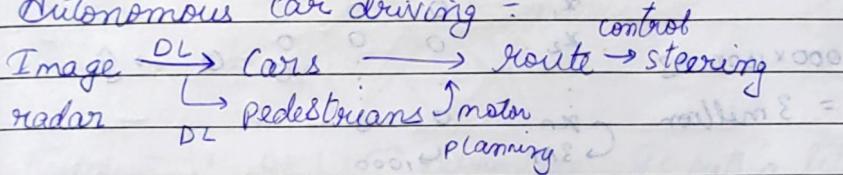
output
 $x \rightarrow y$

- ↳ Let the data speak $x \rightarrow y$ ↳ May need large amount of data
- ↳ Less hand-designing of components needed ↳ Excludes potentially useful hand-designed components

⇒ Applying end-to-end deep learning:

Key question: Do you have sufficient data to learn a functⁿ of the complexity needed to map x to y ?

⇒ Autonomous car driving:



↳ Use DL to learn individual component

↳ Carefully choose $x \rightarrow y$ depends on what tasks you can get data for.

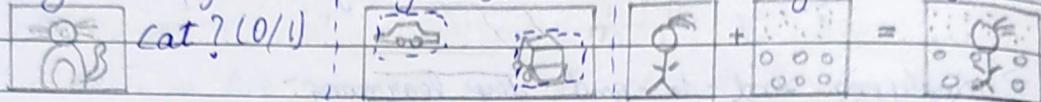
Image ——————> steering

DEEP LEARNING - Course - 4Convolutional Neural Networks

Week 1

- * Foundations of Convolutional Neural Network - Computer Vision :

Image classification | Object Detection | Neural Style Transfer



- ⇒ Deep learning on large images :

Cat? (0/1) i.e. $64 \times 64 \times 3 = 12288$ low res image

$$\begin{array}{c} \text{Cat? (0/1)} \\ \text{1000} \times 1000 \times 3 \\ = 3 \text{ million} \end{array} \quad \begin{array}{l} \text{weights} = (1000, 3M) \\ x_1 \quad 0 \quad 0 \\ x_2 \quad 0 \quad 0 \\ \vdots \quad \vdots \quad \vdots \\ x_n \quad 0 \quad 0 \\ \hookrightarrow 3M \quad \downarrow 1000 \end{array} \quad 0 \longrightarrow g$$

- ⇒ Edge Detection Example :-

Computer Vision Problem :-

* How it works

Vertical
Edges
Horizontal
Edges

- ⇒ Vertical Edge detection :-

$$\begin{array}{c} \text{first} \quad \text{Block} \quad \text{of } 4 \times 4 \quad \text{matrix} \\ \boxed{\begin{array}{|ccc|} \hline 3 & 1 & 2 & 4 \\ 1 & 5 & 8 & 7 \\ 2 & 2 & 1 & 3 \\ 0 & 1 & 3 & 2 \\ \hline \end{array}} \end{array} \quad \begin{array}{c} \text{convolut}^n \\ \downarrow \\ 3 \times 3 \end{array} \quad \begin{array}{c} 1 \ 0 \ -1 \\ 1 \ 0 \ -1 \\ 1 \ 0 \ -1 \end{array} = \begin{array}{c} -5 \ -4 \ 0 \ 8 \\ -10 \ -3 \ 2 \ 2 \\ 0 \ 2 \ -4 \ -2 \\ -3 \ -2 \ -3 \ -16 \end{array}$$

shift the 3×3 block by right

$$\begin{array}{c} 6 \times 6 \\ \xrightarrow{\text{element wise} \times} \begin{bmatrix} 3 & 0 & 1 \\ 1 & 5 & 8 \\ 2 & 7 & 2 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} = \begin{bmatrix} 3 \times 1 & 0 \times 0 & 1 \times 1 \\ 1 \times 1 & 5 \times 0 & 8 \times 1 \\ 2 \times 1 & 7 \times 0 & 2 \times -1 \end{bmatrix} = \begin{array}{c} 3 \ 0 \ -1 \\ 1 \ 0 \ -8 \\ 2 \ 0 \ -2 \end{array} \end{array} \quad * \text{Sum} \quad \begin{array}{c} 3 \ 0 \ -1 \\ 1 \ 0 \ -8 \\ 2 \ 0 \ -2 \end{array} = -5 \rightarrow 4 \times 4$$

sum of all terms

python : conv-forward

tensorflow : tf.nn.conv2d

Keras : Conv2D

⇒ Vertical Edge selection : ∇ thus detects vertical edges

→ This detects horizontal edges

$$6 \times 6 \text{ or } 2 \times 3 \times 3 \times f \times f \quad n-f+1 \times n-f+1$$

6 x 6 ~~9 x 9~~

3x3 fxf

$$n = f + \frac{1}{4} \times n - f$$

lighter band detects the vertical separation of original image
Strong vertical edge

10	10	10	10	0	0	0	
10	10	10	10	0	0	0	
10	10	10	10	0	0	0	
0	0	0	0	10	10	10	
0	0	0	0	10	10	10	
0	0	0	0	10	10	10	
0	0	0	0	10	10	10	

→

1	1	1					

*

0	0	0					

=

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

→

More lighter parts

more darker parts

\Rightarrow Filters: $\begin{matrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{matrix}$ \rightarrow Advantage: It puts a little bit more weight to the central row, the central pixel and it makes it maybe a little bit more robust.

3	0	-3
0	0	-10
3	0	-2

With these filter we use it for edge detection. But computer has a better understanding of selecting the 9 numbers and it often more robust than hand-picked 9 numbers (i.e 3×3 matrix).

Schaum filter - picked 9 numbers (i.e 3×3 matrix)

18

two layers

Padding :- why we need to used Padding i.e adding extra column of blank pixels :-

Now with
padding 6x6

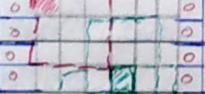
image has
changed

~~50 8 x 8~~

Pad it with

P = padding

$$P=1$$

 i) Shrinking output our picture is originally of 6×6 pixels but shrinks to 4×4 and this may continue as we train high Deep layer neural net.

3×3 ii) Throwing away info from edge : As we only throw away through first cell so we may lack info from that pixel. Now the final Image size will be $n+2p-f+1 \times n+2p-f+1$

0	0	0	0	0	0
0					0
0					0
0					0
0					0
0					0
0					0
0					0
0					0
0	0	0	0	0	0

8 x 8

810

on CP through

1.

~~Sexel, Nov~~

— 1 —

卷之三

⇒ Valid and Same convolutions :-

"Valid" → No padding : $n \times n \star f \times f \rightarrow n-f+1 \times n-f+1$

$$6 \times 6 \star 3 \times 3 \rightarrow 4 \times 4$$

"Same" → Pad so that output size is the same as the input,

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

$$\therefore \text{for "same" padding } n+2p-f+1 = n$$

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

↳ In CV f is almost odd
as it has a central pixel

$$3 \times 3, 5 \times 5, 7 \times 7 \dots$$

⇒ Strided convolution :- [padding p , stride s \Rightarrow i.e. take s jumps as $s=2$]

$$\begin{array}{c} \text{Take 2 steps} \\ \text{1st} \quad 2 \quad 3 \\ \boxed{1 \quad 2 \quad 3} \end{array} \star \begin{array}{c} \text{1st} \quad 2 \quad 3 \\ 1 \quad 2 \quad 3 \\ X \end{array} = \begin{array}{c} \lceil \frac{n+2p-f+1}{s} \rceil \times \lceil \frac{n+2p-f+1}{s} \rceil \\ \frac{7+0-3+1}{2} = 3 \times 3 \\ \text{floor}(7.1) = 7 \end{array}$$

$7 \times 7 \quad 3 \times 3 \quad 3 \times 3$
 $n \times n \quad f \times f \quad s$
(say $s=2$)

⇒ Convolution of RGB image :-

* convolution for a 3D-image RGB images :

$$\begin{array}{ccc} \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} & \star & \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \\ \text{height} \quad \downarrow \quad \text{width} \quad \text{channel} & \xrightarrow{\text{Have to be equal}} & 4 \times 4 \end{array}$$

⇒ Multiple filters :-

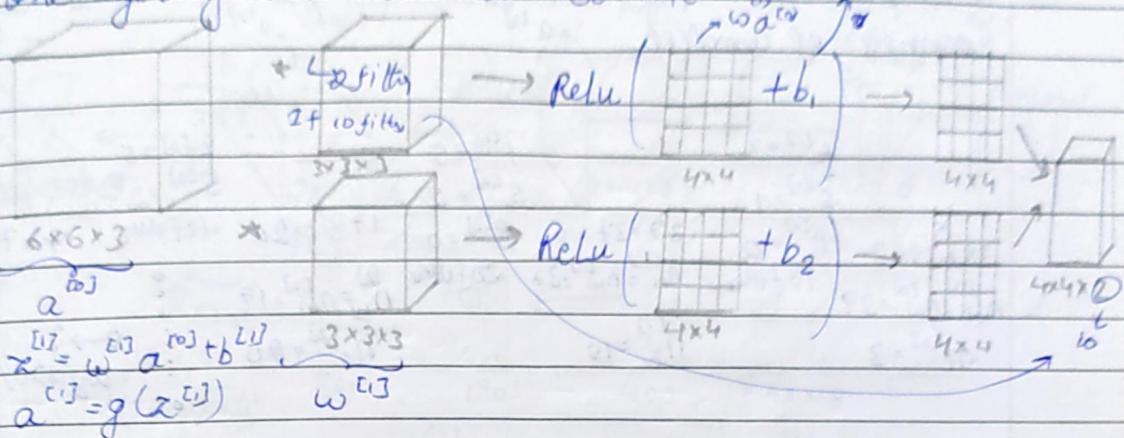
$$\begin{array}{c} \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \quad \star \quad \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \\ 6 \times 6 \times 3 \quad f \times f \times n_c \quad 3 \times 3 \times 3 \quad 4 \times 4 \quad n-f+1 \times n-f+1 \times n_c \end{array}$$

vertical
edge

Horizontal
edge

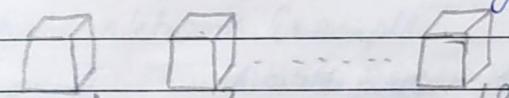
no. of
filters

⇒ One layer of a convolutional Network \dagger



⇒ Number of parameters in one layer \ddagger :

Q) If you have 10 filters that are $3 \times 3 \times 3$ in one layer of a neural network, how many parameters does that layer have?



$3 \times 3 \times 3$ i.e. 27 parameters + bias = 28

$$\therefore 28 \times 10 = 280 \text{ parameters}$$

⇒ Summary of notation \ddagger :

If layer l is a convolution layer:

$$f^{[l]} = \text{filter size} \quad \text{Input} = n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$$

$p^{[l]}$ = padding

$s^{[l]}$ = stride

$$\text{Output} = n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]} \quad \boxed{n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]} + 1}{s^{[l]}} \right\rfloor}$$

Each filter is:

$$f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$$

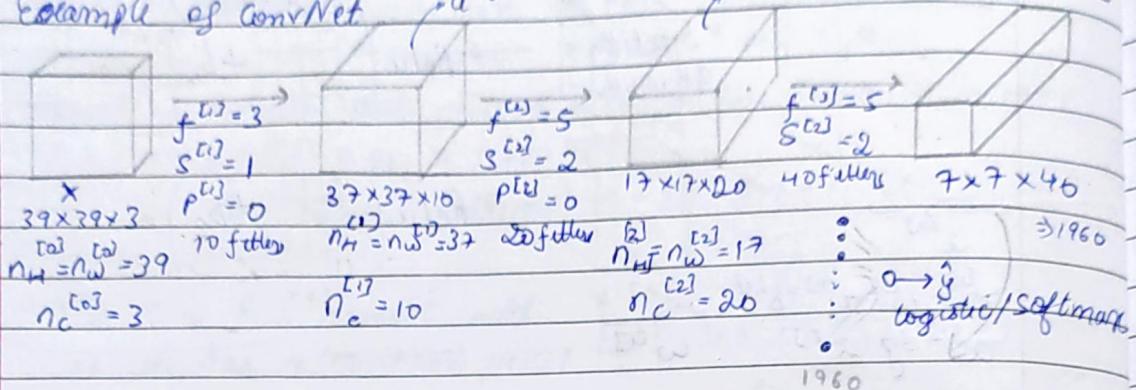
$$\text{Activations} : a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]} \quad | \quad A^{[l]} = m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

Weights : $f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]} \rightarrow n_C^{[l]}$ no. of filters in layer l .

$$\text{Bias} : n_C^{[l]} \rightarrow (1, 1, 1, n_C^{[l]})$$

\Rightarrow A simple convolution network example : Formula $\rightarrow \frac{n+2p-f+1}{s}$

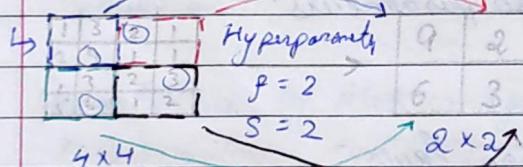
Example of ConvNet



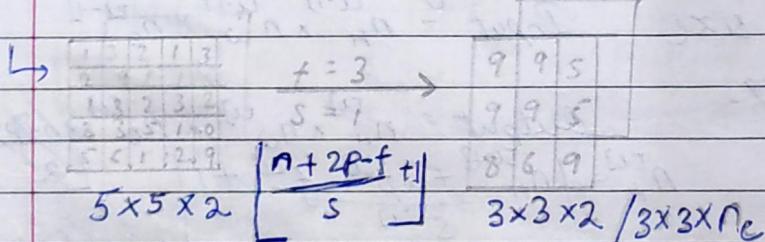
\Rightarrow Types of layers in a convolutional network :

- \hookrightarrow Convolution (conv)
- \hookrightarrow Pooling (POOL)
- \hookrightarrow Fully Connected (FC)

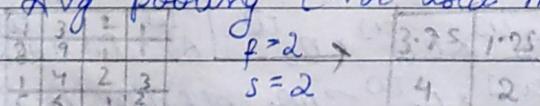
\Rightarrow Pooling layers : Max Pooling



Take maximum number from each 2×2 bracket.



\hookrightarrow Avg pooling (Not used much)



\Rightarrow Summary of pooling : Hyperparameters $n_H \times n_W \times n_C$

$\hookrightarrow f \rightarrow$ filter size ; Say $f=2, 3=2 / f=3, s=2$

$\hookrightarrow s \rightarrow$ stride

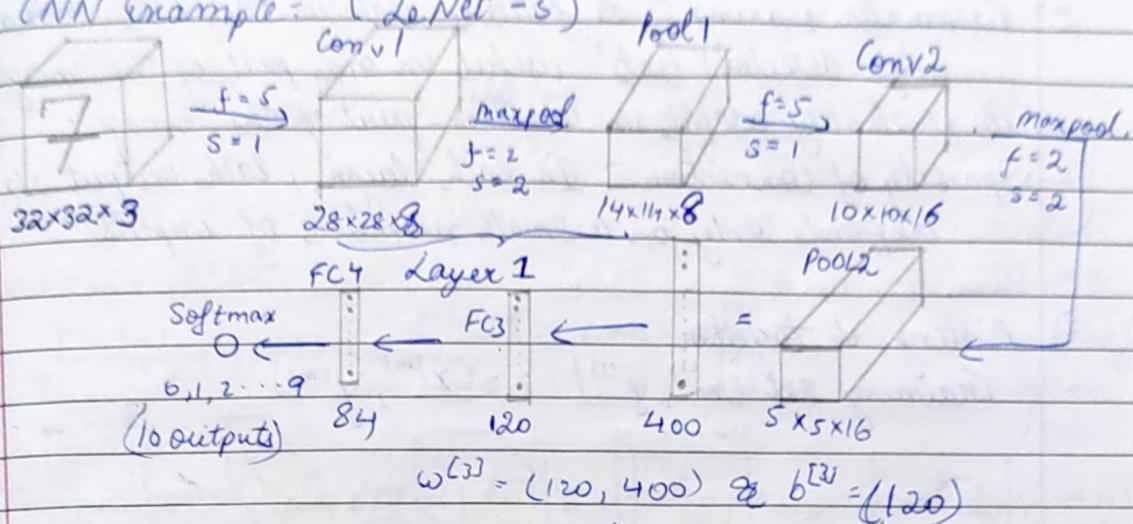
$\hookrightarrow p \rightarrow$ Padding, ($p=0$) after $\left[\frac{n_H-f+1}{s} \right] \times \left[\frac{n_W-d}{s} + 1 \right]$

\hookrightarrow Max or average pooling

No parameters to learn !! (Fixed functⁿ)

$\times n_C$

⇒ CNN Example : (LeNet - 5)



As we go deeper in the neural network n_w & n_b will decrease
& Number of channels will increase;

⇒ Neural Network Example :

	Activation Shape	Activation Size	# parameters
Input	$(32, 32, 3)$	$3072 \text{ a}^{[0]}$	0
Conv1 ($f=5, s=1$)	$(28, 28, 8)$	6272	608
Pool L1	$(14, 14, 8)$	1568	0
Conv2 ($f=5, s=1$)	$(10, 10, 16)$	1600	3216
Pool L2	$(5, 5, 16)$	400	0
FC3	$(120, 1)$	120	48120
FC4	$(84, 1)$	84	10164
Softmax	$(10, 1)$	10	850

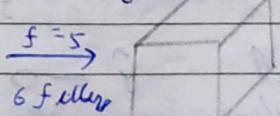
⇒ Why convolutions : why convolution layers have advantage over fully connected layers (FC)

i) Parameter Sharing

ii) Sparsity of connection

So if we create a neural network of
○ ○ $3072 \times 4704 \approx 14M$ it is huge

$$\frac{32 \times 32 \times 3}{3072}$$



$$\frac{28 \times 28 \times 6}{4704}$$

$$3072$$

$$3 \times 5 \times 5 + 1 = 26$$

$$\text{filter weight} = 26$$

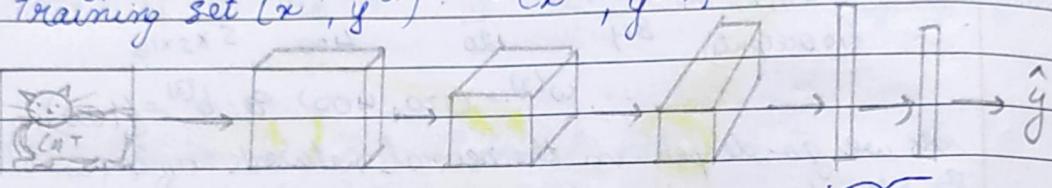
$$6 \times 26 = 156 \text{ parameters}$$

So we can see why conv has few features

- ↳ Parameter sharing - A feature detector (such as vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.
- ↳ Sparsity of connection - In each layer, each output value depends only on a small number of inputs.

- ⇒ Putting it together:

Training set $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$



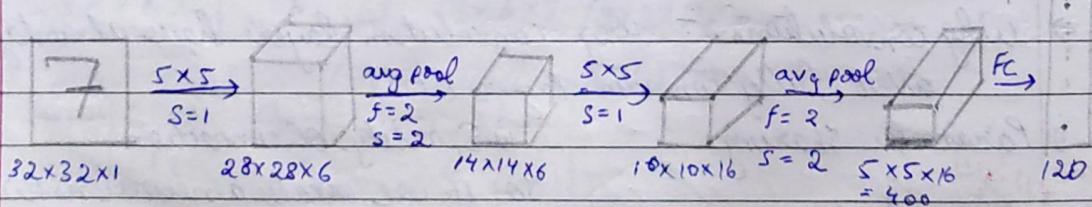
$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

use gradient descent to optimize parameters to reduce J .

Week - 2

- * Case Studies = Outline
- ⇒ Classic Networks = $\begin{matrix} \rightarrow \text{Res Net}(152) \\ | \\ \rightarrow \text{Inception} \\ | \\ \rightarrow \text{LeNet-5} \\ | \\ \rightarrow \text{Alex Net} \\ | \\ \rightarrow \text{VGG} \end{matrix}$

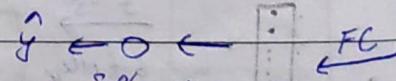
⇒ LeNet - 5



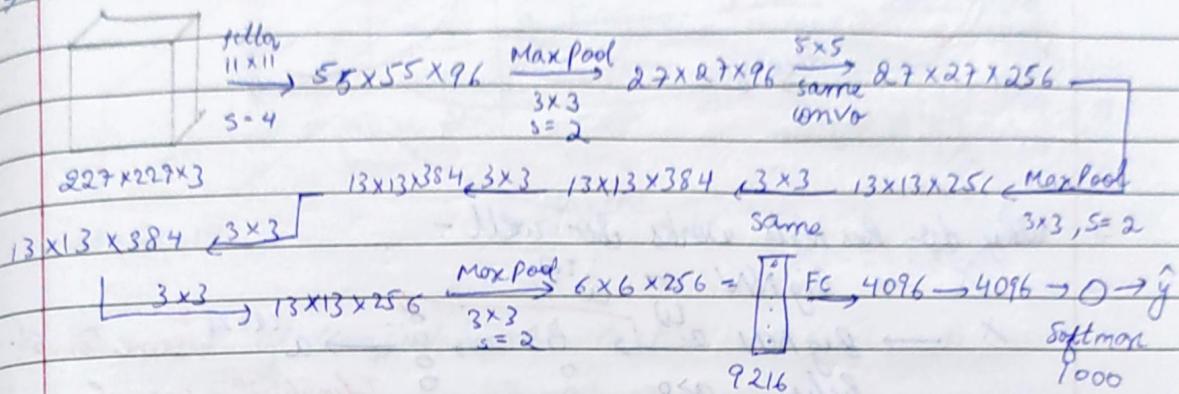
from left to right $(n_H, n_W) \leftarrow g(n)_H$

Conv pool conv pool fc fc output

In early days they used sigmoid/tanh instead of ReLU



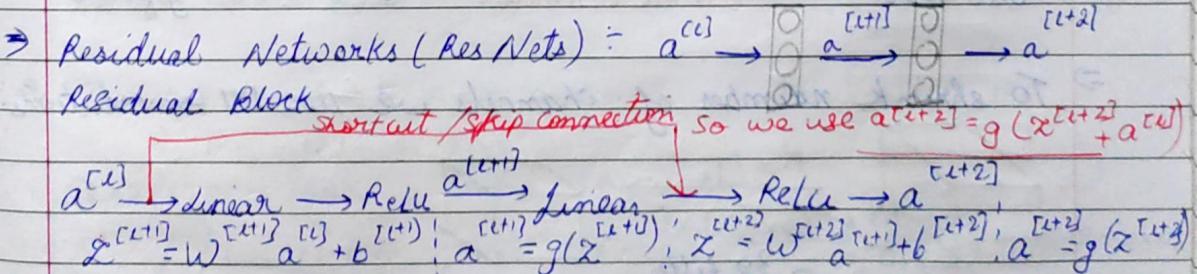
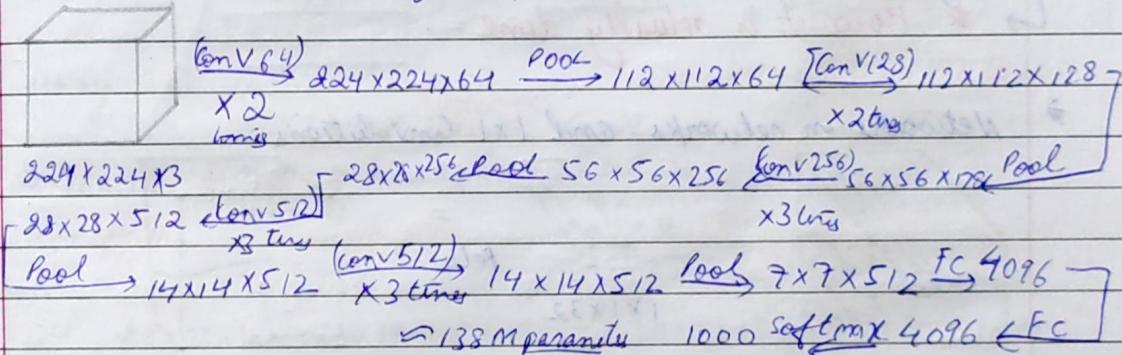
⇒ AlexNet :-



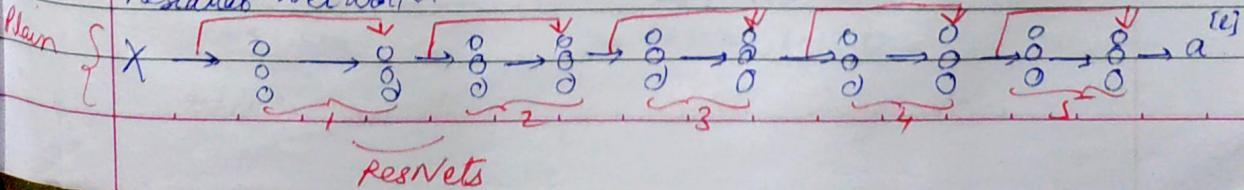
- ⇒ It is similar to LeNet but much bigger.
- ⇒ We use ReLU.
- ⇒ Multiple GPUs are used.
- ⇒ Local Response Normalization (LRN) Not used Nowadays.

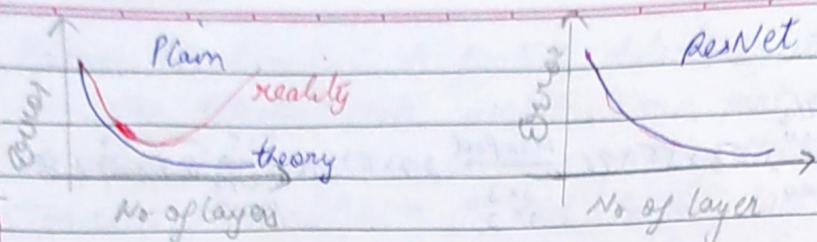
⇒ $\#$ means 16 layers which has weight

⇒ VGG - 16 \div Conv = 3x3 filter, $s=1$, Same Conv & Maxpool = 2x2, $s=2$



⇒ Residual Network:-





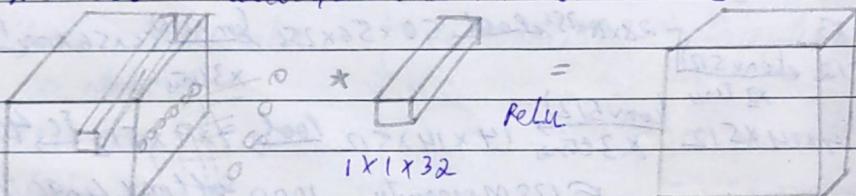
⇒ Why does ResNet work so well?

$$\begin{aligned}
 X &\rightarrow \text{Big NN} \rightarrow a^{[l]} \\
 X &\rightarrow \text{Big NN } a^{[l]} \xrightarrow{\text{ReLU}} a^{[l]} \xrightarrow{\text{Identity } f^l} a^{[l+1]} \xrightarrow{\text{some Conv}} a^{[l+2]} \\
 a^{[l+2]} &= g(a^{[l+2]} + a^{[l]}) \\
 &= g(a^{[l+2]} w^{[l+2]} + b^{[l+2]} + a^{[l]}) \\
 &\hookrightarrow \text{If } w^{[l+2]} = 0 \text{ & } b^{[l+2]} = 0 \quad a^{[l+2]} = g(a^{[l]}) \\
 &= a^{[l]} \text{ [ReLU]}
 \end{aligned}$$

⇒ Why do residual networks work?

↳ How it is actually done

⇒ Networks in networks and 1×1 convolutions:

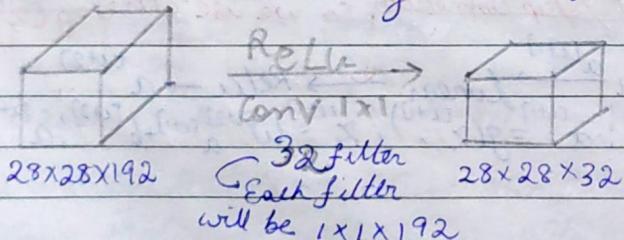


$6 \times 6 \times 32$

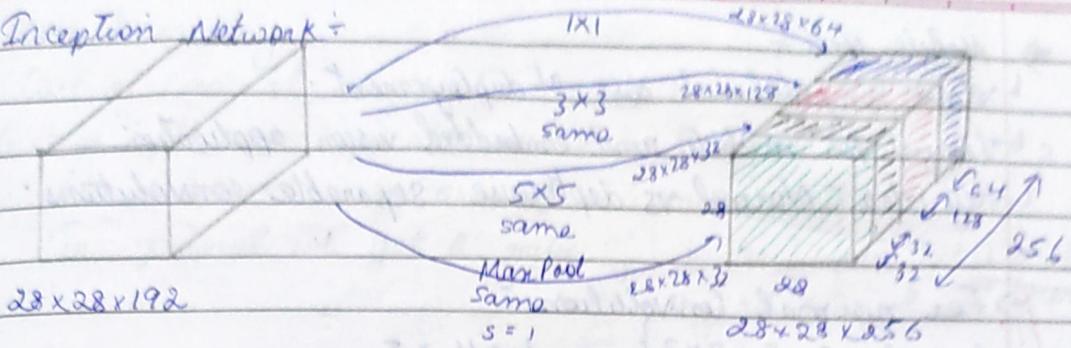
Network in Network

$6 \times 6 \times 32$ No. of filters

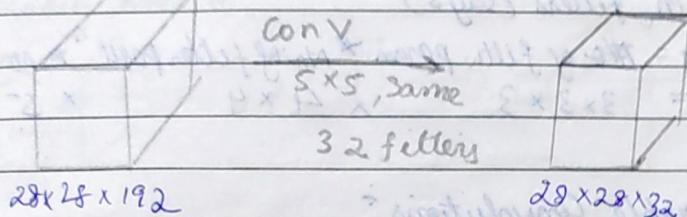
⇒ To shrink number of channels we use 1×1 convolutions:



⇒ Inception Network :



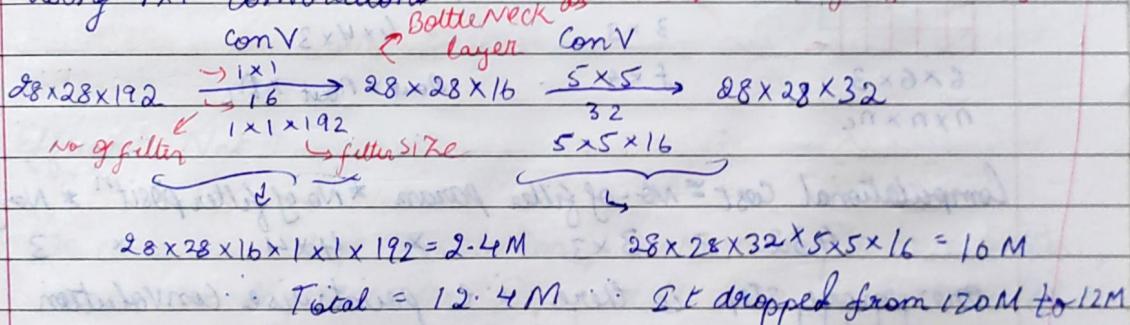
⇒ The prob. of computational cost :



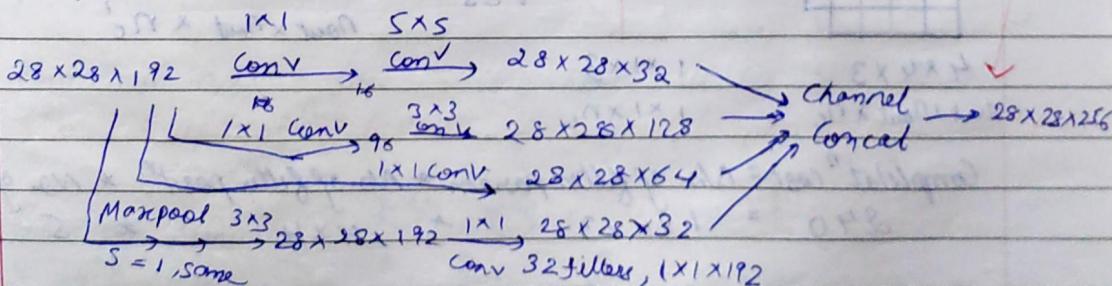
Hence, we have 32 filters of $5 \times 5 \times 192$.
Total $28 \times 28 \times 32 \times 5 \times 5 \times 192 \approx 120M$

So we can see it computationally expensive. We will use 1×1 convolution and it will reduce it by $1/10^{th}$.

⇒ Using 1×1 convolutions :



⇒ Inception module :



* How does it actually looks : Google Net.

⇒ Mobile Net -

↳ Low computational cost at deployment.

↳ Useful for mobile and embedded vision application.

↳ Key idea: Normal vs depthwise - separable convolutions.

⇒ For normal convolutions:

$$\begin{matrix} 6 \times 6 \times 3 & * & 3 \times 3 \times 3 \\ n \times n \times n_c & & n \times n \times n_c \end{matrix} = 4 \times 4 \times 5$$

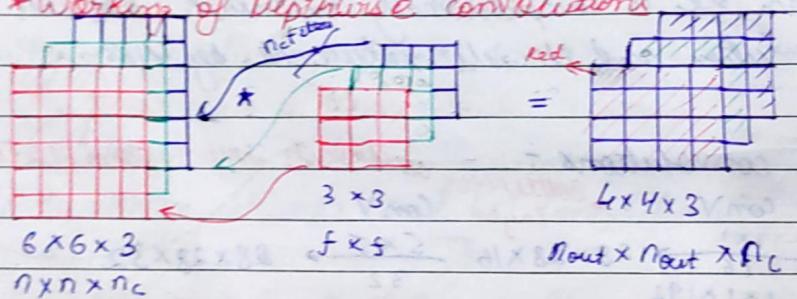
\downarrow
 n_c' filters (say 5)

$$\text{Computational cost} = \# \text{ of filter param} * \text{No of filter posit}^n * \text{no of filters}$$

$$2160 = 3 \times 3 \times 3 \times 4 \times 4 \times 5$$

⇒ i) Depthwise Separable Convolutions:

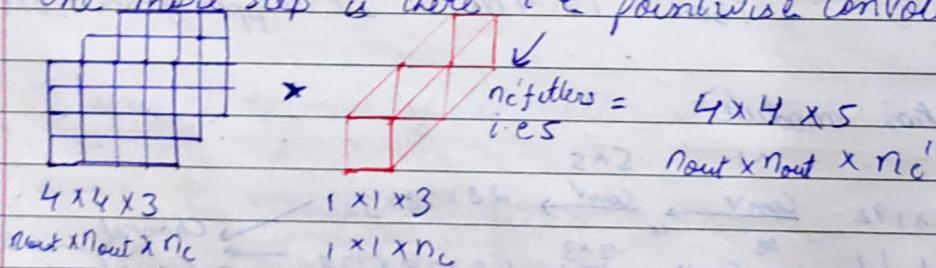
* Working of Depthwise convolution



$$\text{Computational Cost} = \text{No. of filter param} * \text{No of filter posit}^n * \text{No of filters}$$

$$432 = 3 \times 3 \times 4 \times 4 \times 3$$

ii) One more step is there i.e pointwise convolution



$$\text{Computat' cost} = \text{No. of filter param} * \text{No of filter posit}^n * \text{No of filters}$$

$$240 = 1 \times 1 \times 3 * 4 \times 4 * 5$$

⇒ Cost Summary :-

Cost of Normal convolution = 2160

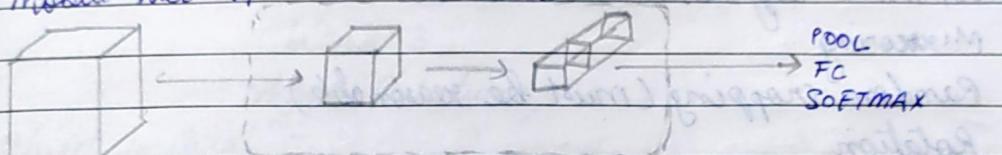
Cost of depth wise separable convolution = $432 + 240 = 672$

∴ Normal convos were around 31% more expensive

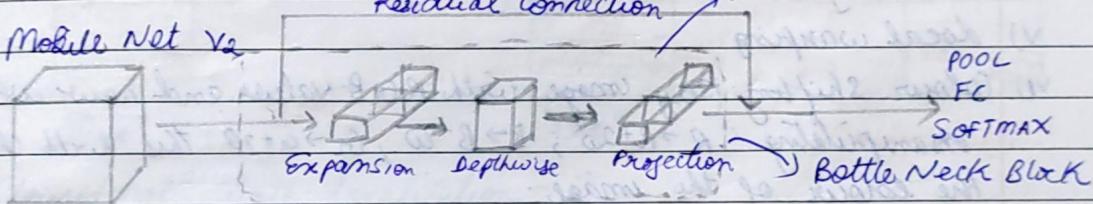
In general we get a ratio of $\frac{1}{n_0} + \frac{1}{f^2}$ i.e. $L + L \approx 31\%$
 $n_0 \quad f^2 \quad 5 \quad 3^2$ (Here)

⇒ Mobile Net Architecture :-

↳ Mobile Net v1



↳ Mobile Net v2



* How does bottleneck work

⇒ Efficient Net :-

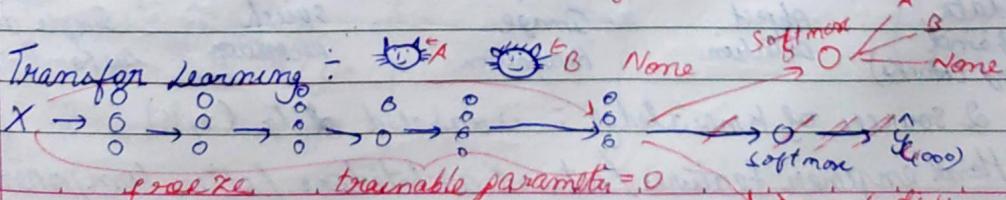
* How Efficient Network works

↳ We can use High Resolution, make a deeper architecture, make the layers wider. So which one is best suited or you can also do compound scaling i.e. simultaneously performing 2 tasks.

* Practical advice for using conv-net :-

* Practical application of ResNets using github open source via CMD

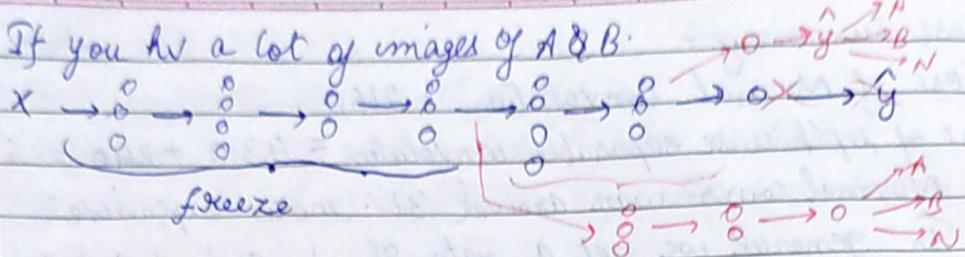
⇒ Transfer Learning :-



If u hv less pics of A & B - freeze = 1

→ save & disk

⇒ If you hv a lot of images of A & B.



If u hv sufficient data u can train your own Net.

⇒ Data augmentation :-

common augmentation method :-

i) Mirroring

ii) Random cropping (must be reasonable)

iii) Rotation

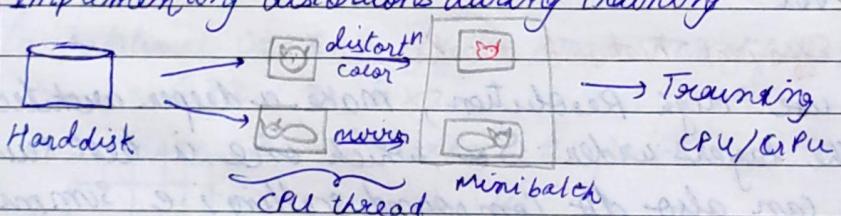
iv) Shearing $\square \rightarrow \square'$

v) local warping

vi) Colour shifting (say image with RGB values and now we manipulating $R \rightarrow R+20$; $B \rightarrow B-20$, $G \rightarrow G+20$ this will change the colour of the image.)

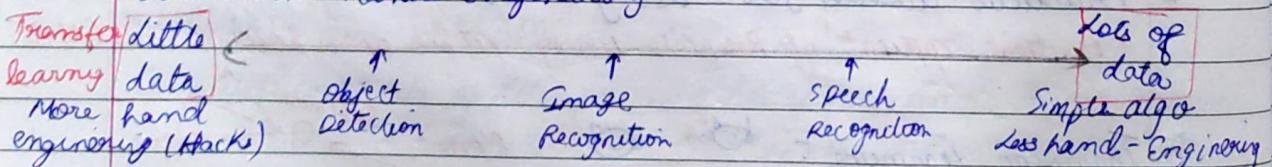
→ Alex Net Paper for colour augmentation "Pct colour augmentation".

⇒ Implementing distortions during training :-



⇒ State of Computer Vision :-

Data v/s hand engineering



2 sources of knowledge :- i) Labeled data (x, y)

ii) Hand engineer features / network architecture / other components .

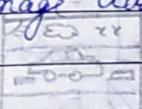
- 3) Tips for doing well on benchmarks/winning competition:
 - i) Ensembleing: Train several networks independently and average their outputs (\hat{y})
 - ii) Multi crop at test time: Run classifier on multiple versions of test images and average results.
 - iii) Use open source code: Use architecture of networks published in the literature / Use open source implementation if possible / Use pretrained models and fine-tune on your dataset.

Week - 3

Object Detection

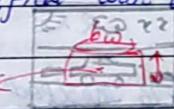
* Object localization:

Image classification



'Car'

Classification with localization

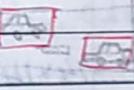


(b_x, b_y, b_w, b_h)

"Car"

$$\begin{aligned} b_x &= 0.5 \\ b_y &= 0.2 \\ b_h &= 0.3 \\ b_w &= 0.4 \end{aligned}$$

Detection



multiple object

Classification with localization:

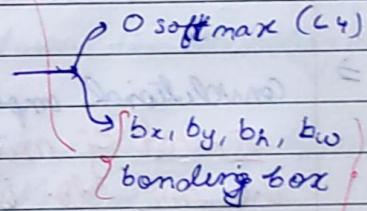
(0,0)

(1,1)

Image



Conv Net



Defining the target label y :

1. pedestrian

Need to output b_x, b_y, b_w, b_h class label (1-4)
 P_c → is there any object

2. car

So far an image with car

3. Motorcycle

with car

4. Background

No car

$$\hookrightarrow \text{loss function: } L(\hat{y}, y)$$

$$\left\{ \begin{array}{l} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_3 - y_3)^2 \quad \text{if } y = 1 \\ (\hat{y}_1 - y_1)^2 \quad \text{if } y = 0 \end{array} \right\}$$

$$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_h \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

$$y = \begin{bmatrix} b_x \\ b_y \\ b_w \\ b_h \\ 0 \end{bmatrix}$$

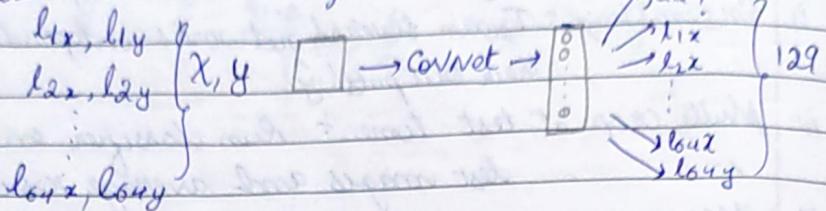
$$y = \begin{bmatrix} ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

No one car

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

No one car

⇒ Landmark Detection : We can find landmarks of nose, eyes, lips of face?



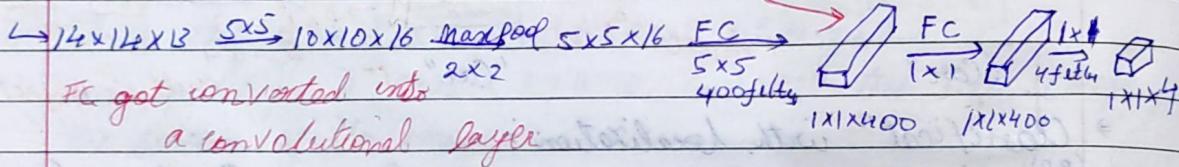
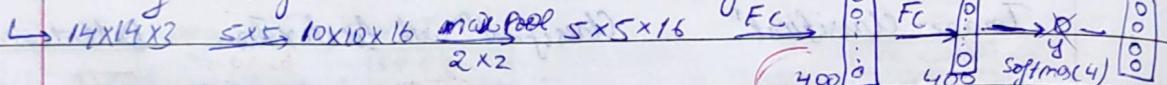
⇒ Object Detection :

Car detection example : we use Sliding window detection

* Sliding window detecting working

⇒ Convolutions involving Sliding windows :

Turning FC layer into convolutional layers :



⇒ "convolutional implement" of Sliding windows :

Instead of using sliding window one by one. What we can do is use those windows altogether and then detect the object

* How it is implemented :

* Real life example

⇒ Object Detection : Bounding box predictions

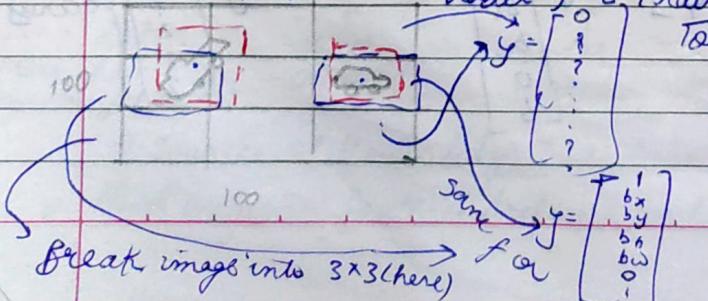
YOLO Algorithm (You only look once)

"." checks mid point

label for training ; for each grid cell :

Pc
b ₁
b ₂
b ₃
b ₄
b ₅
b ₆
b ₇
b ₈

Total output = $3 \times 3 \times 8$

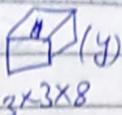


b_1	b_2	b_3
b_4	b_5	b_6
b_7	b_8	b_9
c_1	c_2	c_3

So;



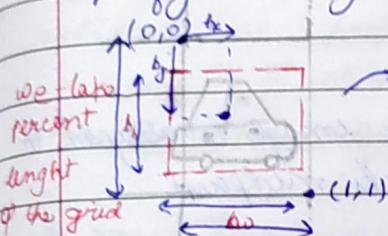
\rightarrow Conv \rightarrow max pool



$3 \times 3 \times 8$

$100 \times 100 \times 3$ So in this algo we get bounding box of the object.
we might use a grid of 19 by 19 instead of 3 by 3.

\Rightarrow Specify Bounding Boxes:



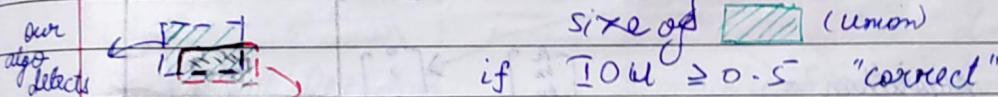
$$y = \begin{bmatrix} b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0.3 \\ 0.4 \\ 0.5 \\ 0.9 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

b1/b0
 0.031
 could
 be greater
 than 1

3×3 grid

\Rightarrow Intersection over union : Evaluating object localization (IoU)

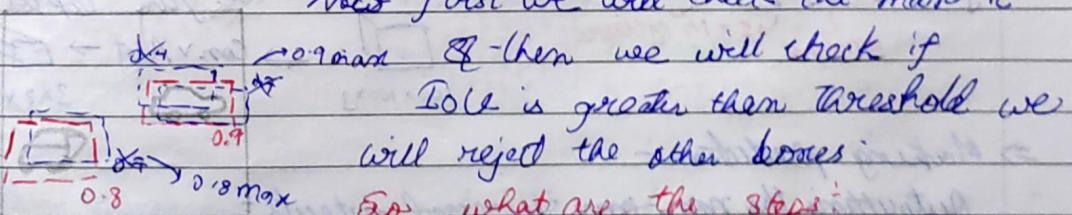
$\text{IoU} = \frac{\text{size of intersection}}{\text{size of union}}$



IoU is a measure of the overlap between two bounding boxes

\Rightarrow Non-max suppression : It makes sure that our objects is detected only once, not multiple times.

Now first we will check the max IoU



Then we will check if IoU is greater than threshold we

will reject the other boxes

So what are the steps:

- Finally queued all boxes with $fc \leq 0.6$
- While there are any remaining boxes
 - Pick out the box with the largest fc ; Output as the prediction
 - Discard any remaining box with $\text{IoU} \geq 0.5$ with the box output in the previous step

⇒ Anchor Boxes : overlapping object.

So we can create 2 anchor box 1 & anchor box 2.



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

y_1 anchor box 1



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

anchor
box 2

↳ Previously - Each object in training image is assigned to grid cell that contains that object's midpoint

Output $y = 3 \times 3 \times 8$

↳ with two anchor boxes - Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IOU (gridcell, anchorbox)

↳ Output $y = 3 \times 3 \times 16$ or $3 \times 3 \times 2 \times 8$

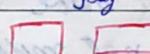
⇒ YOLO Algorithm :

(Training)

1. man

2. car

3. motorcycle



$$y = 3 \times 3 \times 2 \times 8$$

anchors $\rightarrow 3 + \text{No. of classes}$

$$y = \begin{bmatrix} p_x \\ b_x \\ b_y \\ b_w \\ b_h \\ c_1 \\ c_2 \\ c_3 \\ p_c \end{bmatrix} \rightarrow \begin{bmatrix} ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

$$y_1 = \begin{bmatrix} ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

$$y_2 = \begin{bmatrix} ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

what we give input
ConvNet \rightarrow $3 \times 3 \times 16$

⇒ Making prediction :-

Outputting the non-max suppressed outputs :-

↳ For each grid cell, get 2 predicted bounding boxes.

↳ get rid of low probability predictions.

↳ For each class (pedestrian, car., motorcycle) use non-max suppression to generate final prediction.

⇒ Region proposal = R-CNN

* There are many boxes where nothing is to be found, so a R-CNN segmentation algorithm is used and then objects are predicted.

Faster algorithms -

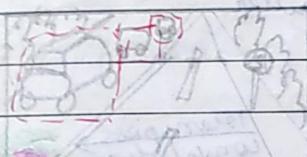
⇒ R-CNN : Propose regions. Classify proposed regions one at a time. Output label + bounding box.

⇒ Fast R-CNN : propose regions. Use convolution implementation of sliding windows to classify all the proposed regions.

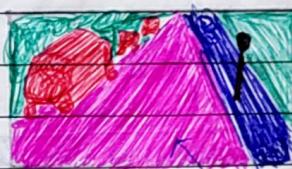
⇒ Faster R-CNN : Use convolutional network to propose regions.

⇒ Semantic Segmentation with U-Net:

Object detection v/s Semantic Segmentation.



Object Detection



For self driving cars instead of detecting road it is better to

Semantic Segmentation find out and every pixel what it represents.

* Other applications like chest X-ray & Brain MRI (Say road)

⇒ Per pixel class label =

3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3
2	2	2	2	2	2	2	2
2	2	2	2	2	2	2	2
2	2	2	2	2	2	2	2
2	2	2	2	2	2	2	2

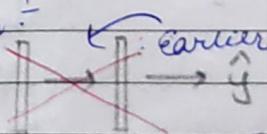
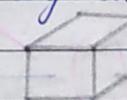
Segmentation Map 1. Car

3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3
2	2	2	2	2	2	2	2
2	2	2	2	2	2	2	2
2	2	2	2	2	2	2	2
2	2	2	2	2	2	2	2

2. Road

3. Building

⇒ Deep learning for Semantic Segmentation =



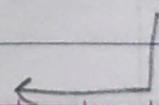
Image

But now we have increased

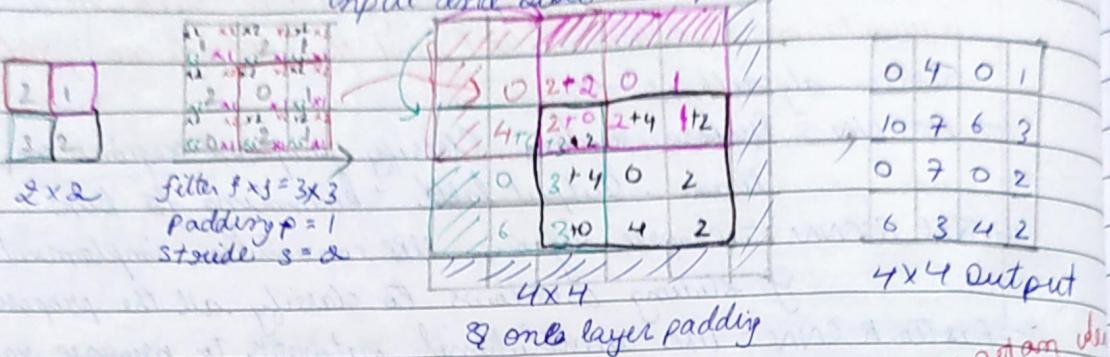
last layers to get our U-Net and semantic

segmentation. For that we use transpose

segmentation which we will study ahead.



=> Transpose Convolution : Basically how do you take 2×2 input and blow it up to 4×4 input.



⇒ U-Net Architecture intuition

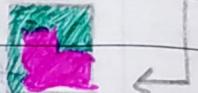
Deep learning for Semantic Segmentation

[large image](#)

In this we can get an idea of what is some cat like image present in our memory.

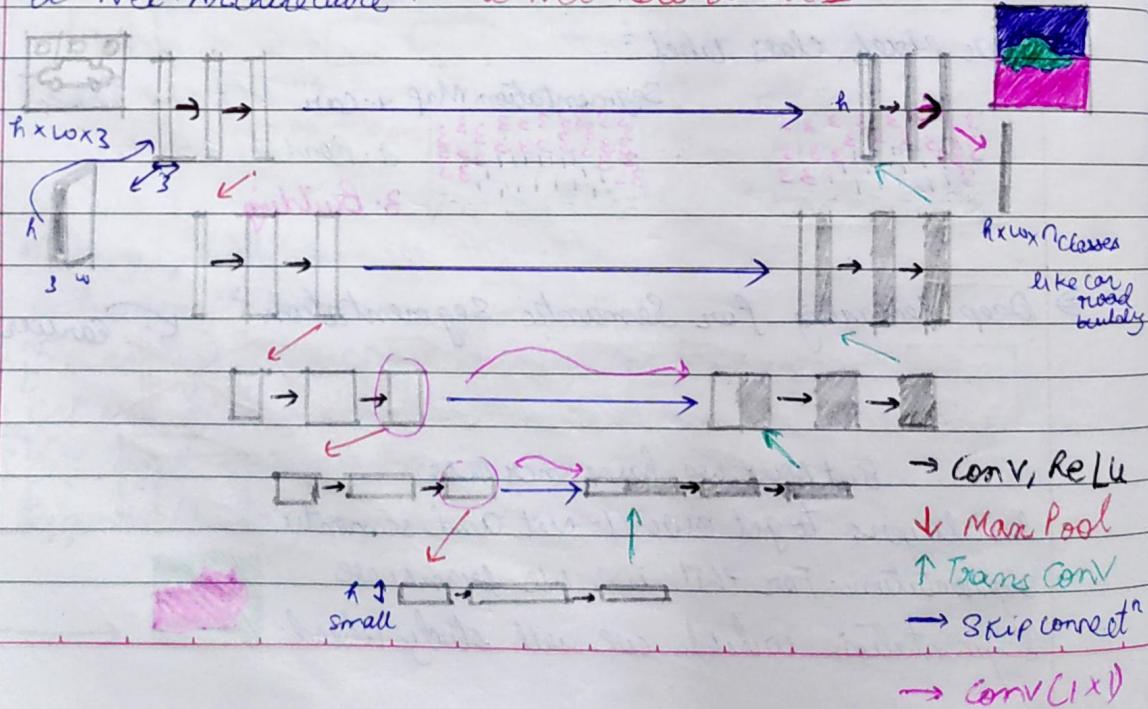
↳ Normal convolutions as size decreases for the first part

Transpose convolution



↳ In smaller one we have lots of spatial information

⇒ U-Net Architecture - A U-Net How it works



Week-4

Special Applications: Face recognition & Neutral Style Transfer

* Face Recognition :

* Face recognition system : Brute system face recognition

Face Verification 1:1

Face recognition 1:K

→ Input image, name / ID

→ Has a database of K persons

→ Output whatever the input image

→ Get an input image

is that of the claimed person

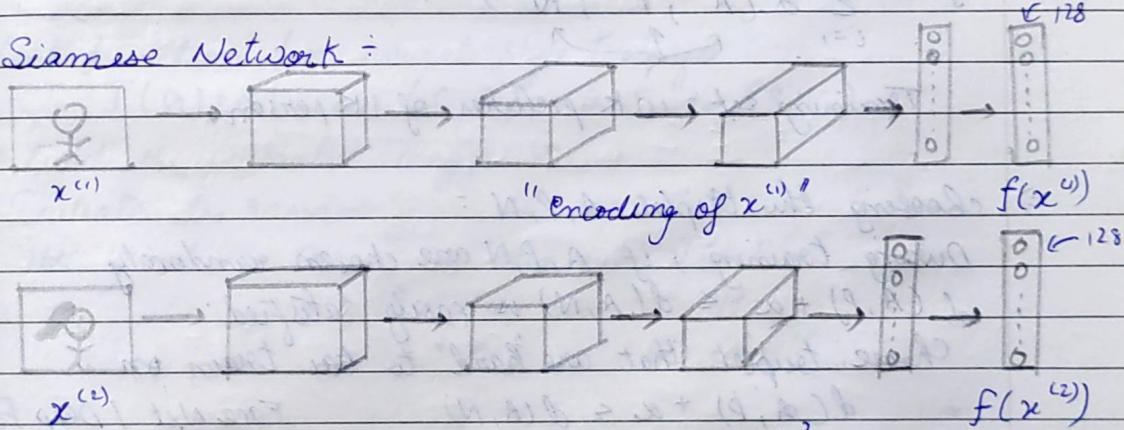
→ Output ID if the image is any of the K persons (or "not recognized").

⇒ One shot learning : we often have less data for recognit/verified
So, we use ab concept of "similarity" function

$d(\text{img 1}, \text{img 2})$ = degree of difference between images

If $d(\text{image 1}, \text{image 2}) \leq \tau$ "Same" } Verification
 $> \tau$ "different" }

⇒ Siamese Network :



$$d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_2^2$$

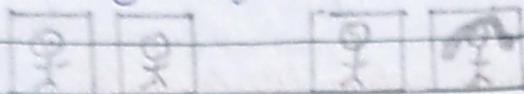
⇒ Goal of learning :

Parameters of NN define an encoding $f(x^{(i)})$ ↗ 128
Learn parameters so that :-

→ If $x^{(i)}, x^{(j)}$ are the same person $\|f(x^{(i)}) - f(x^{(j)})\|_2^2$ is small

→ If $x^{(i)}, x^{(j)}$ are the different person $\|f(x^{(i)}) - f(x^{(j)})\|_2^2$ is large

⇒ Face Recognition = Triplet Loss → It uses multi 3 images → Anchor, Positive, Negative
Learning objective



Anchor Positive Anchor Negative
(A) (P) (A) (N)

(say α)

We want $\|f(A) - f(P)\|^2 + \alpha \leq \|f(A) - f(N)\|^2$ Not big enough
 $d(A, P) = \alpha$ $d(A, N) = \alpha / \sqrt{2}$ we keep alpha to prevent $0 - 0 \leq 0$
 $\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 \leq 0 - \alpha$
So our net work will give reasonable output $\xrightarrow{\text{Reasonable}}$

⇒ loss function :

Given 3 images A, P, N :

$$L(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)})$$

Training set : 10K pictures of 1K person

⇒ choosing the triplets A, P, N :

During training, if A, P, N are chosen randomly

$L(A, P) + \alpha \leq L(A, N)$ is easily satisfied

Choose triplet that are "hard" to train on

$$d(A, P) + \alpha \leq d(A, N)$$

Face Net / Deep Face

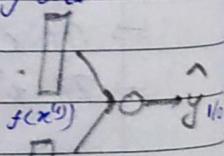
$$\frac{d(A, P)}{\downarrow} \approx \frac{d(A, N)}{\uparrow}$$

⇒ Face verification and Binary classification :

Using a Siamese network of 2 person and finally add a single dense layer for binary classification

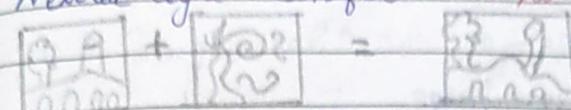
$$y = \sigma \left(\sum_{k=1}^n w_k [f(x^{(i)})_k - f(x^{(j)})_k] + b \right)$$

$$\frac{f(x^{(i)})_k - f(x^{(j)})_k}{f(x^{(i)})_k + f(x^{(j)})_k}^2 = X^2$$



$$f(x^{(i)})$$

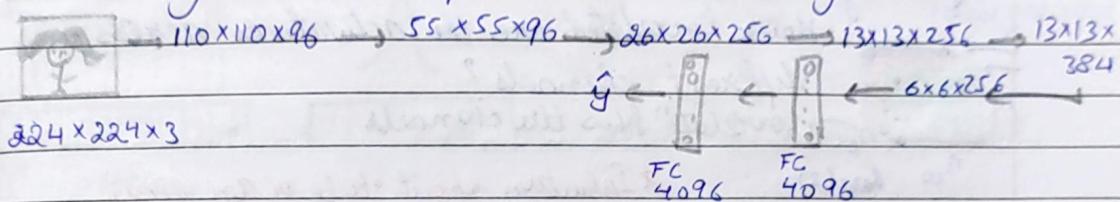
* Neural Style Transfer: * How neural style Transfer works



Content (C) Style (S) Generated Image (G)

⇒ What are deep ConvNet learning?

Visualizing what a deep Network is learning:



Pick a unit in layer 1. Find the nine image patches that maximize the unit's activation. Repeat for other units

* How to visualize deep layers

⇒ Cost function for Neural Style Transfer:

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

→ Find the generated image G :

i) Initiate G randomly; $G : 100 \times 100 \times 3$

ii) Use gradient descent to minimize $J(G)$

$$G := G - \frac{\partial}{\partial G} J(G) \quad * \text{How it works}$$

⇒ Content cost function:

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

→ Say you use hidden layer l to compute content cost.

→ Use pre-trained ConvNet (Eg: VGG Network)

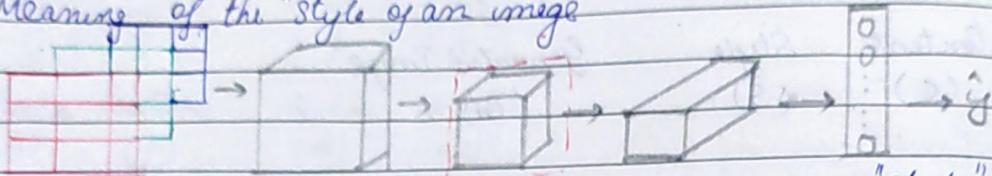
→ Let $a^{[l](C)}$ & $a^{[l](G)}$ be the activation of layer l on the images.

→ If $a^{[l](C)}$ & $a^{[l](G)}$ are similar both images have similar content.

$$J_{\text{content}}(G, G) = \frac{1}{2} \|a^{[l](C)} - a^{[l](G)}\|^2$$

⇒ Style cost function:

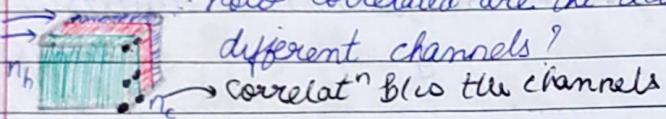
Meaning of the "style" of an image



say you are using layer l 's activation to measure "style".

Define style as correlation between activations across channels.

How correlated are the activations across



now here (2)
it can be $5, 7, 8 \dots$

* Intuition about style of an image

⇒ Style matrix:

Let $a_{i,j,k}^{[l]}$ = activation at (i, j, k) . $G_l^{[l]}$ is $n_c^{[l]} \times n_c^{[l]}$
 $\therefore G_{K,K'}^{[l]}$, \rightarrow It compares how related images are in
 channel $K \& K'$ $K = 1 \dots n_c^{[l]}$

$$G_{K,K'}^{[l](S)} = \sum_{i=1}^{n_c^{[l]}} \sum_{j=1}^{n_c^{[l]}} a_{ijk}^{[l](S)} \cdot a_{ijk'}^{[l](S)} \quad \left\{ \begin{array}{l} \text{"Gram matrix"} \\ (G_l) \end{array} \right.$$

$$G_{K,K'}^{[l](G)} = \sum_{i=1}^{n_c^{[l]}} \sum_{j=1}^{n_c^{[l]}} a_{ijk}^{[l](G)} \cdot a_{ijk'}^{[l](G)}$$

$$\beta J_{\text{style}}^{[l]}(S, G_l) = \| G_l^{[l](S)} - G_l^{[l](G)} \|_F^2 \quad \text{--- Frobenius norm}$$

does not matter much

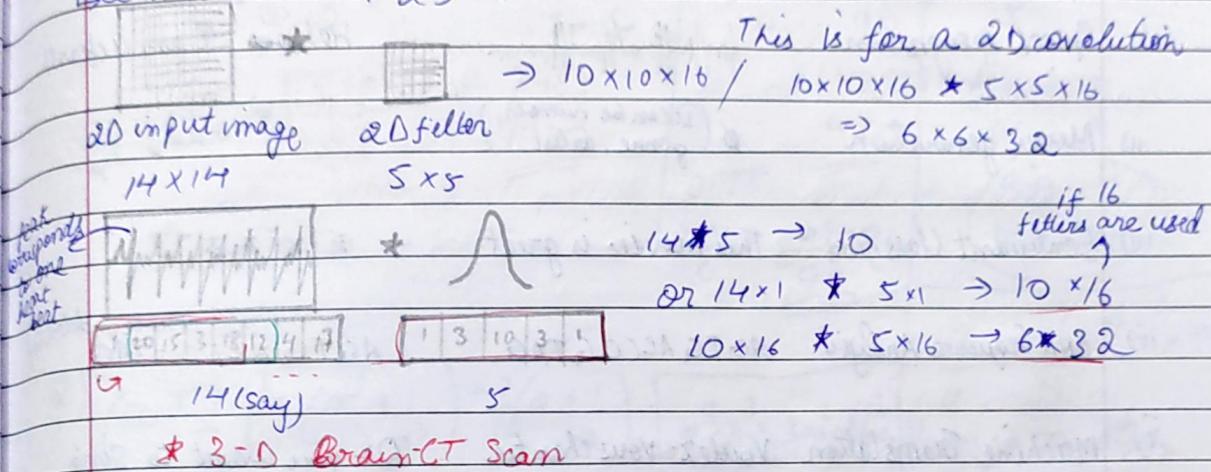
$$= \frac{1}{n_c^{[l]} \times n_c^{[l]}} \sum_{k=1}^{n_c^{[l]}} \sum_{k'=1}^{n_c^{[l]}} (G_{K,K'}^{[l](S)} - G_{K,K'}^{[l](G)})^2$$

$$J_{\text{style}}(S, G_l) = \sum_i \lambda^{[l]} J_{\text{style}}^{[l]}(S, G_l)$$

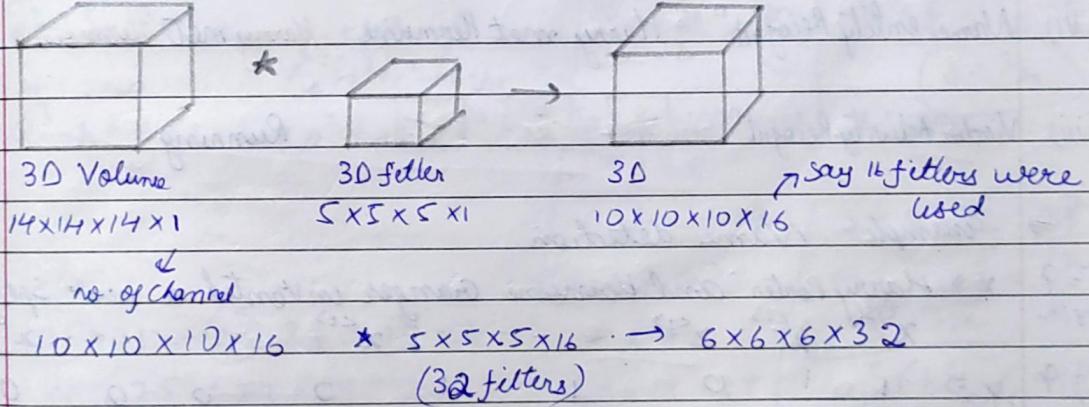
$$\underbrace{J(G_l)}_{C_l} = \alpha J_{\text{content}}(G_l, G_l) + \beta J_{\text{style}}(S, G_l)$$

⇒ 1D and 3D Convolutional Networks
1D & 3D generalization of models

↳ Convolutions in 2D & 1D:



↳ 3D convolution:



Brain CT, MRI Scan, Movies slices of it comes under 3D convolution