

Convolutional Neural Networks

Computer vision

Edge detection

More edge detection.

Use of edge detection at

We use many 3x3 filter to detect

Edge detection picture

edge

Picture matrix * matrix \Rightarrow new matrix

through this matrix machine learn where
the pixels are located and learned
the patterns.

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

left dark bottom light right dark

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

Horizontal

horizontal
edge

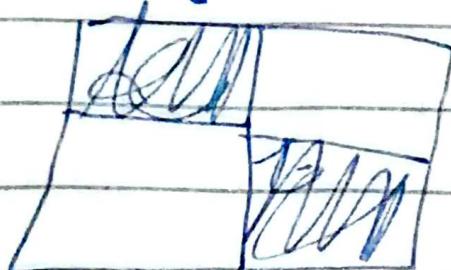
line detected

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

y

\therefore	$\begin{bmatrix} 10 & 10 & 10 & 10 \\ 10 & 10 & 10 & 10 \\ 10 & 10 & 10 & 10 \\ 10 & 10 & 10 & 10 \end{bmatrix}$
--------------	------------------------------------------------------------------------------------------------------------------

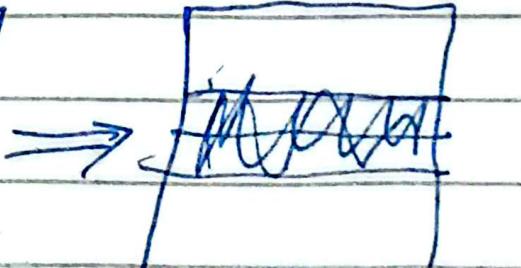
X



edge we get in pixels

$$Y \times X =$$

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



Other filters

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

$$\begin{bmatrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$$

Sobel filter

Scharr filter

(

vertical dark

Sobel

Basic edge detector
feature extraction

Scharr

Same as sobel.
Better at edge detection
regardless of direction

Scharr is more
accurate.

The Padding

When we mask or filter
the size of picture shrinks and
sometimes we don't want to lose them

* Cause only center things are
detected not on the border.

Padding

we get the same size of image where borders are detected.

$$6 \times 6 \times 3 \xrightarrow{\text{padding}} \text{Pad } 6 \times 6$$

How much pad: padding

Valid convolution: no filter: $6 \times 6 \times 3 \Rightarrow 4 \times 4$

Same convolution & output size will be same in input and output size.

$\Rightarrow P=2$ will be required to get full Picture.

T# Strided convolution

* Reduce output size :-

feature map shrink, but with learning weights.

Reduce computation

Compute larger patterns fast.

Working

stride = 2

filters jump jump 2 columns in every row.

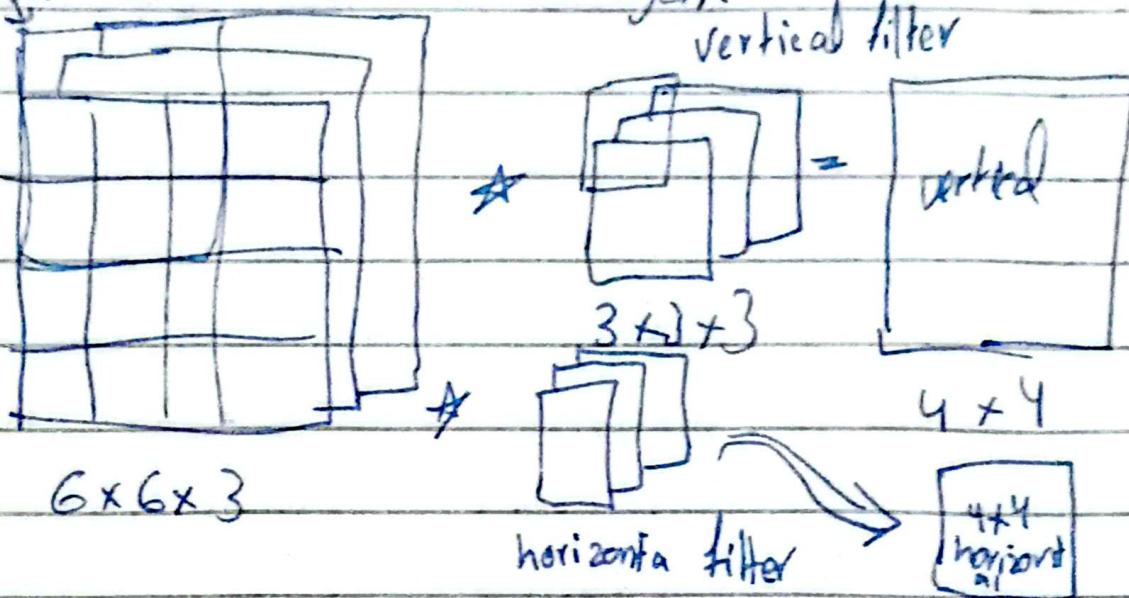
By stroke

$$7 \times 7 + 3 \times 3 = 3 \times 3$$

$$\left[\frac{n+2p-f}{s} + 1 \right] \times \left[\frac{n+2p-f-1}{s} \right]$$

Convolution over RGB image

27 boxes at once on all three layers.



work as normal filter but compute
all three RGB layer at once and add them
then move to next rows or columns

One layer of a convolutional network

If layer $[l]$ is a convolutional layer,

$f^{[l]}$ = filter size

$p^{[l]}$ = padding

$s^{[l]}$ = stride

If $n^{(l)}$ = number of filters.

Each filter is =

Activation:

Weights:

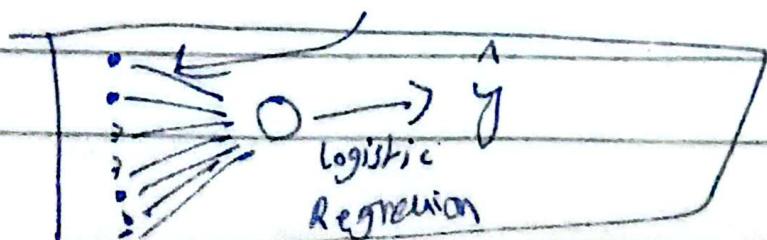
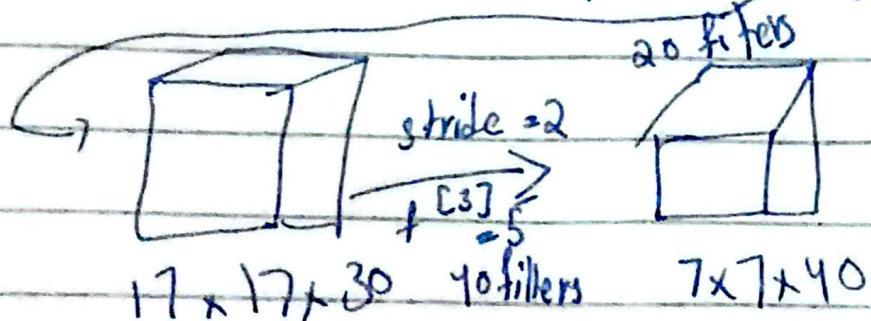
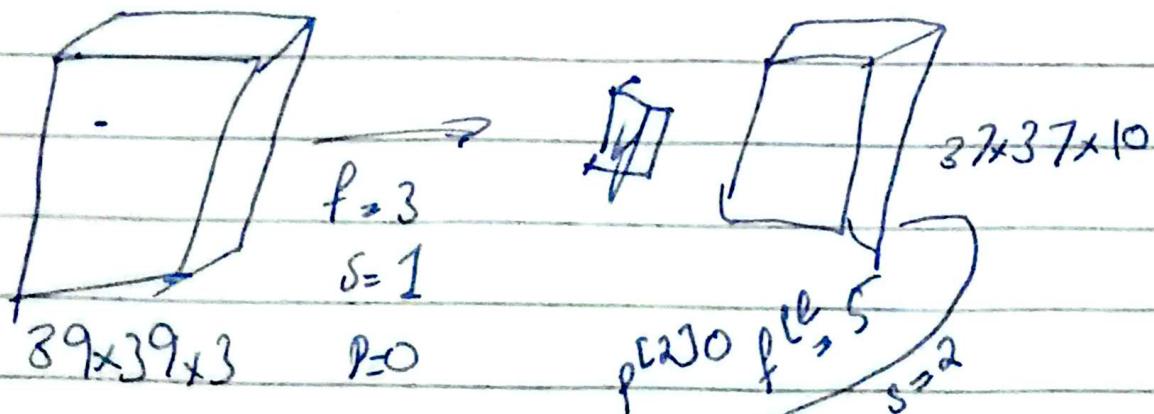
bias:

Simple Convolutional Network example:-

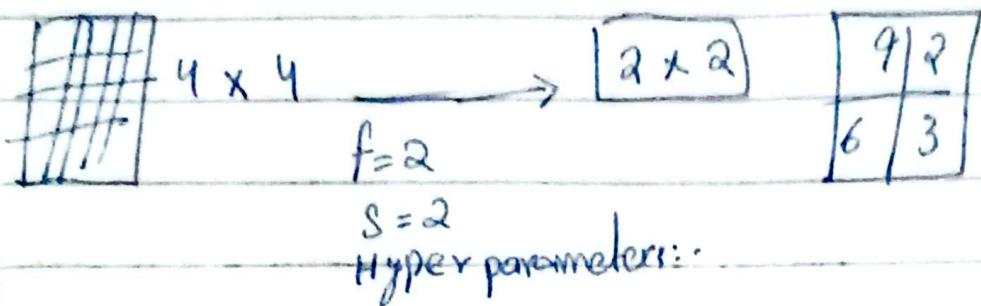
$p=0$ mean valid convolutions.

Conv Net

$39 \times 39 \times 3$

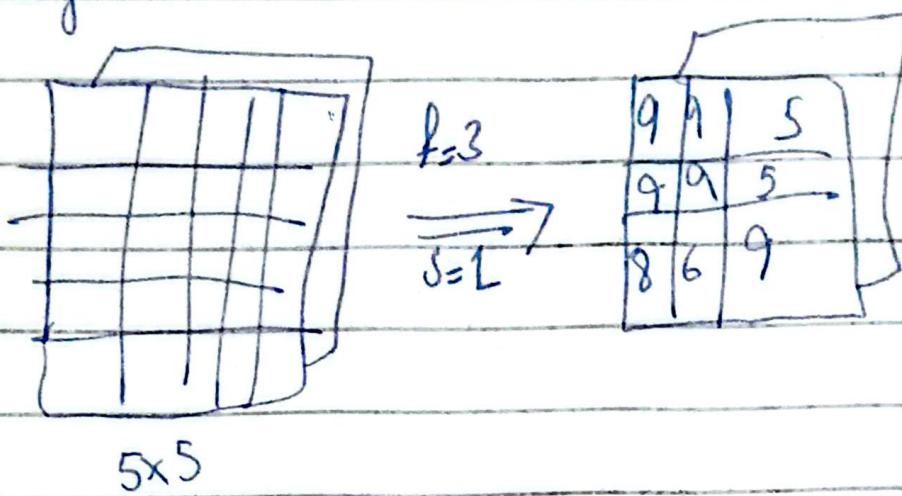


Pooling Layer: Max pooling



It extract features and reduce size of image.

e.g.



5x5

This takes maximum number in a filter.

Type of Average filter

this take average in a filter

Result:

Max pooling is used more than the average pooling

Summary

- Convolution layer: learn filters or to detect edge filter, corners.

- Pooling layer:-

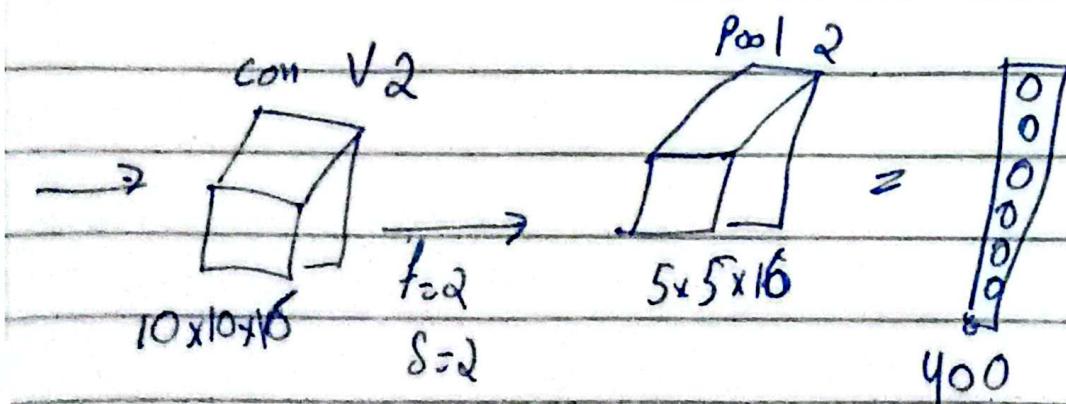
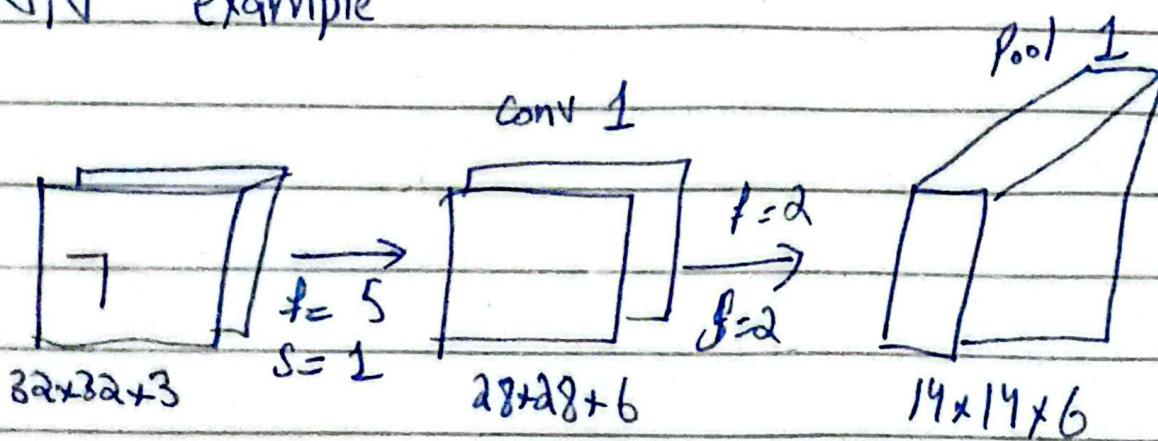
Makes images smaller while keep important info - Reduces computation

Fully connected layer:-

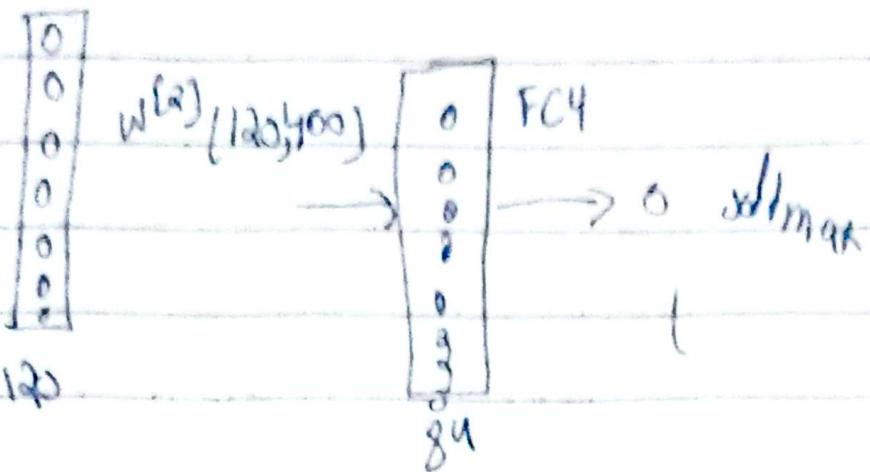
works on normal NN

use activation function features to predict.

CNN example

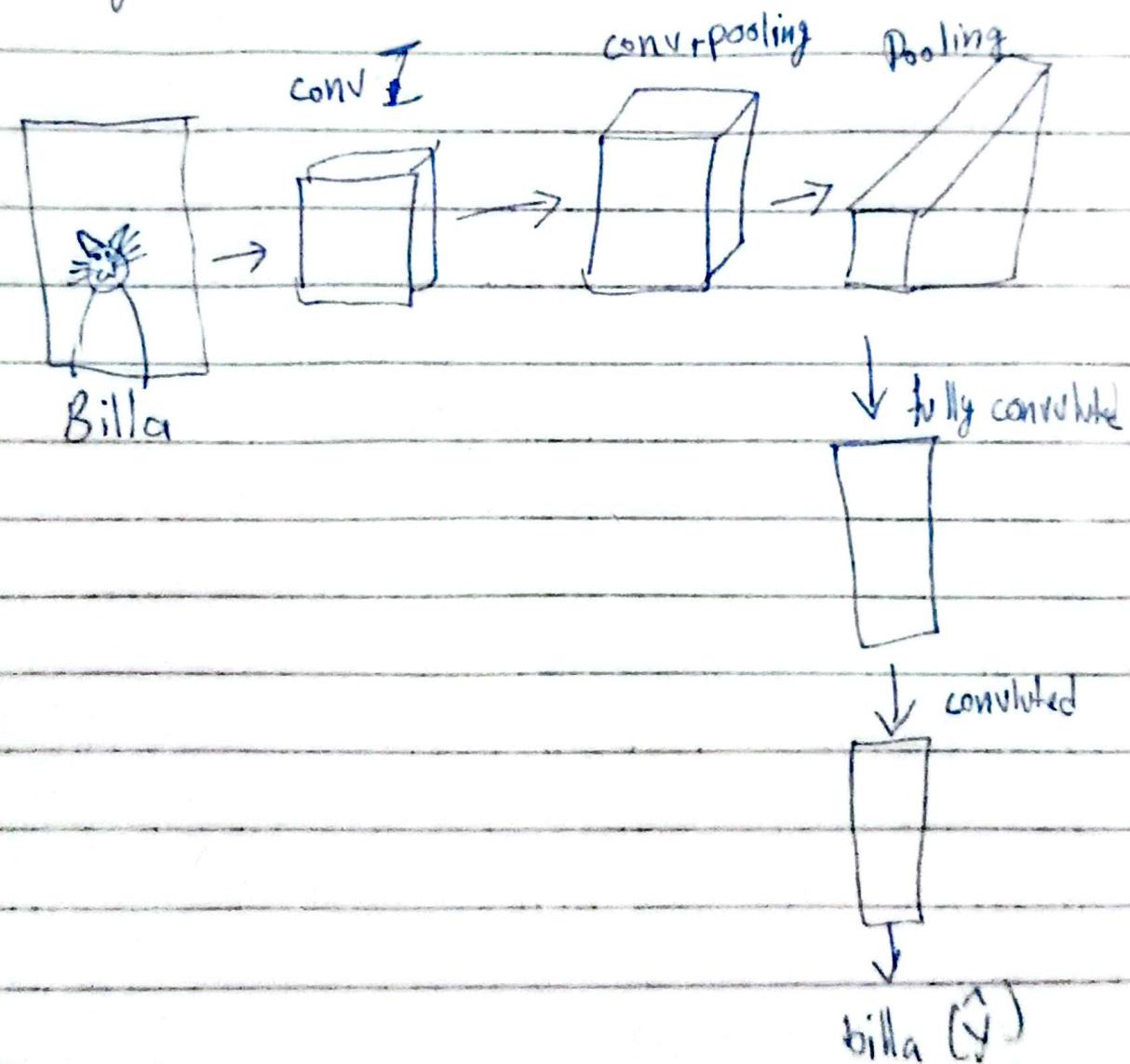


FC3

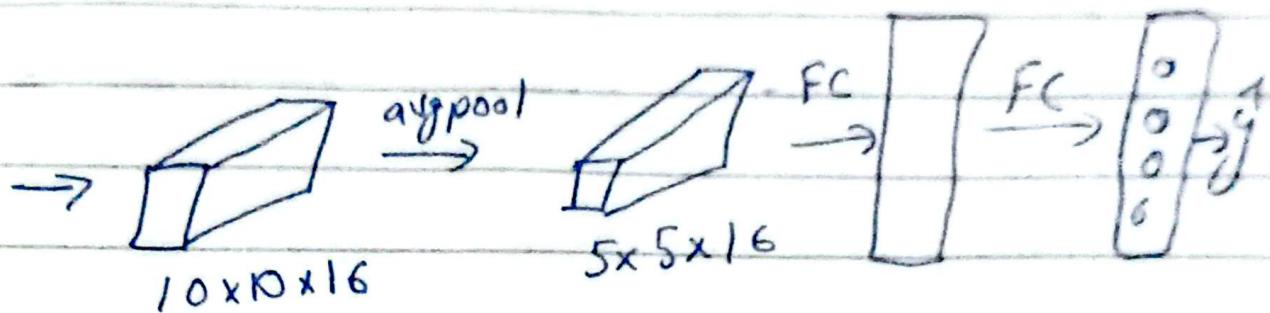
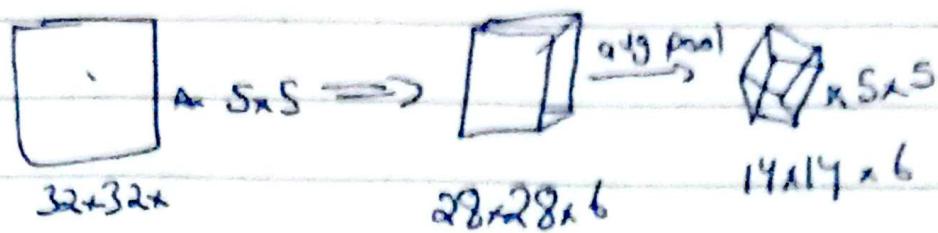


Research:-

Q1 Why convolutions:-



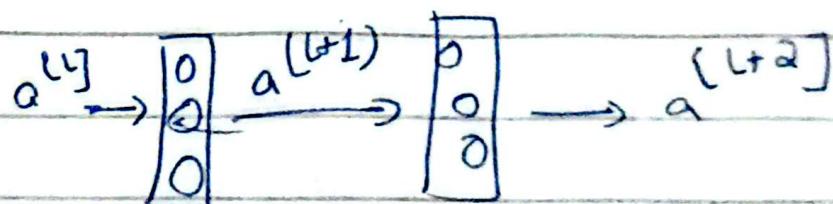
Old Net (Le Net) 1998



VGG16 (16 layer)

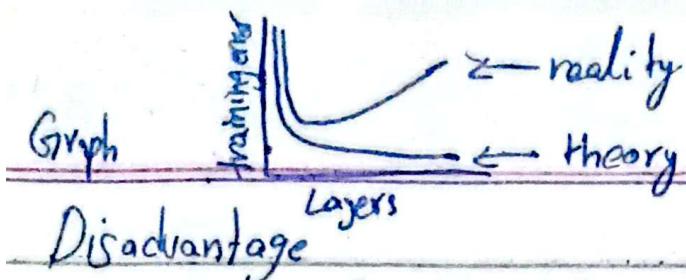
Alex Net

Resnet



$a^{(l)}$ $\xrightarrow{\text{linear}}$ Linear \rightarrow Relu $\xrightarrow{a^{(l+1)}}$ Linear \rightarrow Relu $\xrightarrow{a^{(l+2)}}$

Resnet allow us to train much deeper

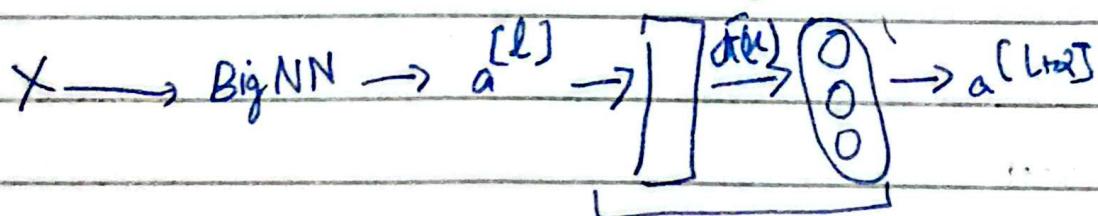


Disadvantage

- 1) High Computational Cost
- 2) Overfitting on small dataset
- 3) Only Extremely deep versions give only marginal accuracy gain while increasing training cost.
- 4) Deeper model are slower during prediction

Why Resnets Work?

This skips connection



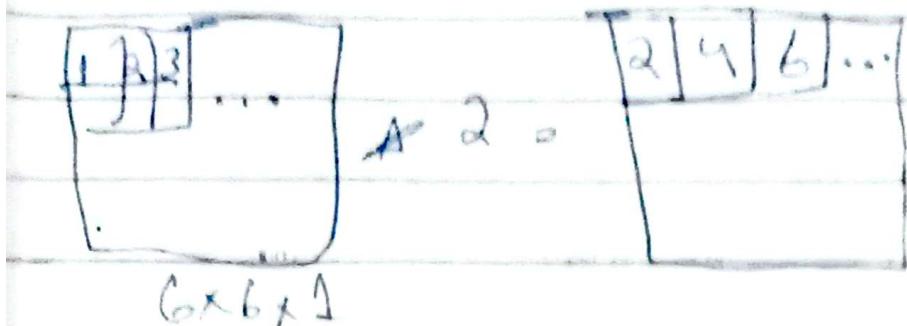
skipped layer

give more accuracy

not hurting performance

They are skip layer bcz the skip over one or more layers and is added directly

1x1 Convolution and Network in Network



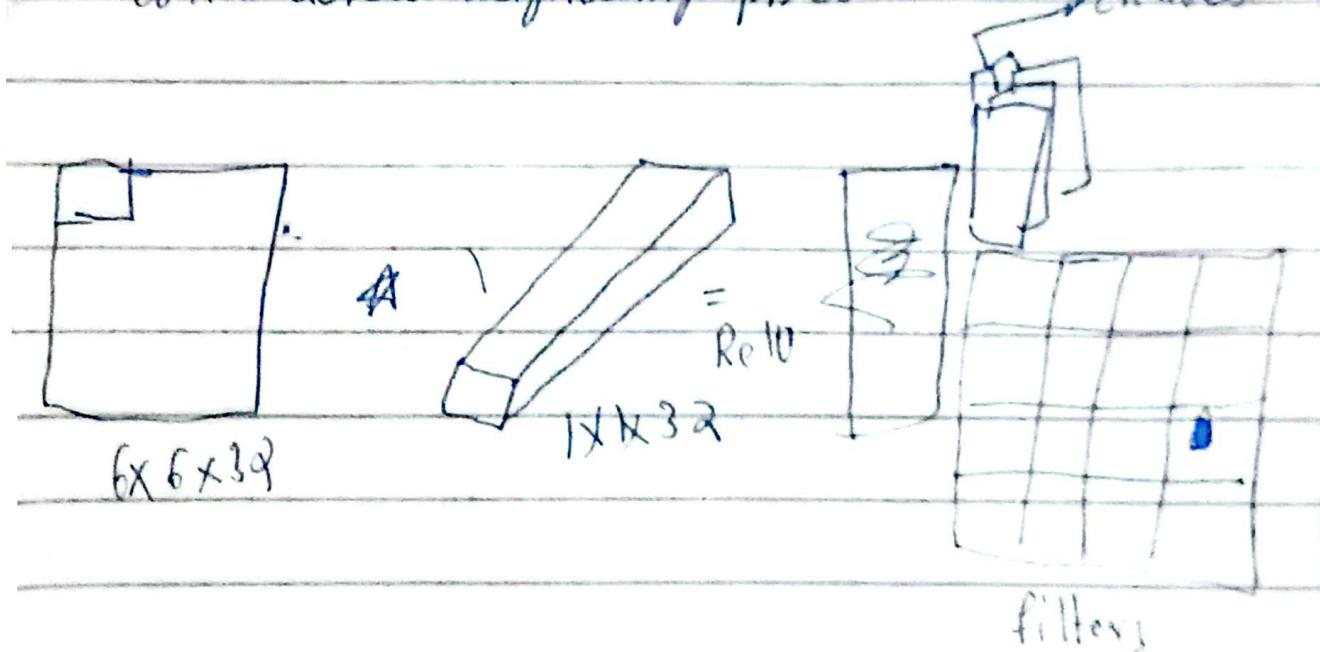
1x1 Convolution

* Filter size 1 height x 1 width

* It slides over the input feature map one spatial position at a time.

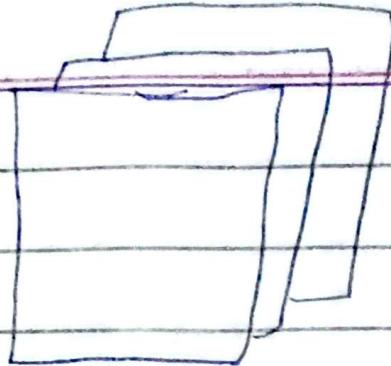
* It only mixes information across channels,

* don't across neighboring pixels. → channels



channels = layers

$$6 \times 6 \times 32 \star 1 \times 1 \times 32 =$$

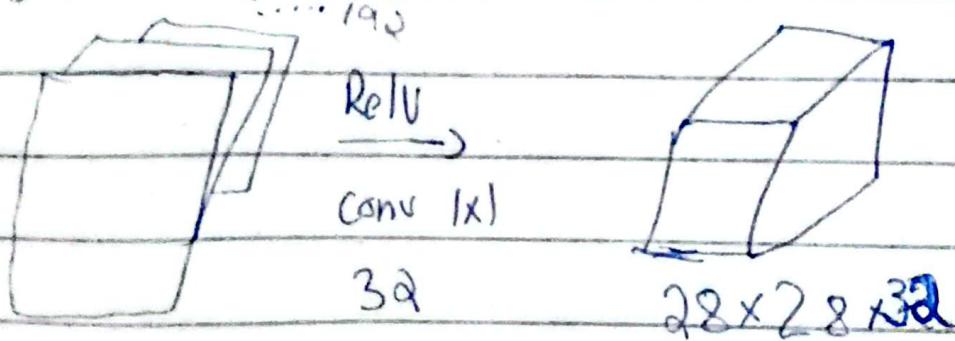


$6 \times 6 \times 1$, filter

one create many
feature map

Like an example

We want to create $1 \times 1 \times 32$
3d filters from 192



$28 \times 28 \times 192$

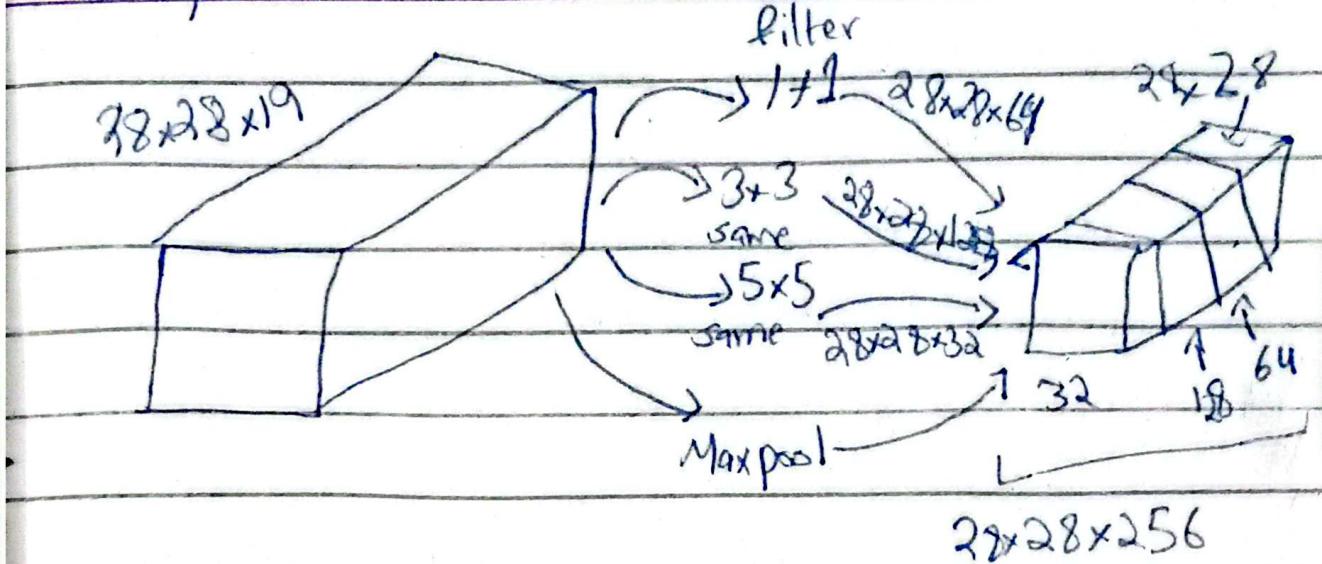
we created

32 filter

features

or channels

Inception Network Motivation



Motivation Before 2014

→ Small filters capture fine details but miss larger context.

→ Larger filters capture more context but are expensive and prone to overfit.

Idea:-

Why not use multiple filter and network decide which is useful.

Why we use it:-

→ To reduce computation cost with 1x1 or bottlenecks.

→ To make network deeper and wider without exploding computation.

3) How it works.

1x1 captures fine details, reduces dimension.

2) 1x1 conv \rightarrow 3x3 conv

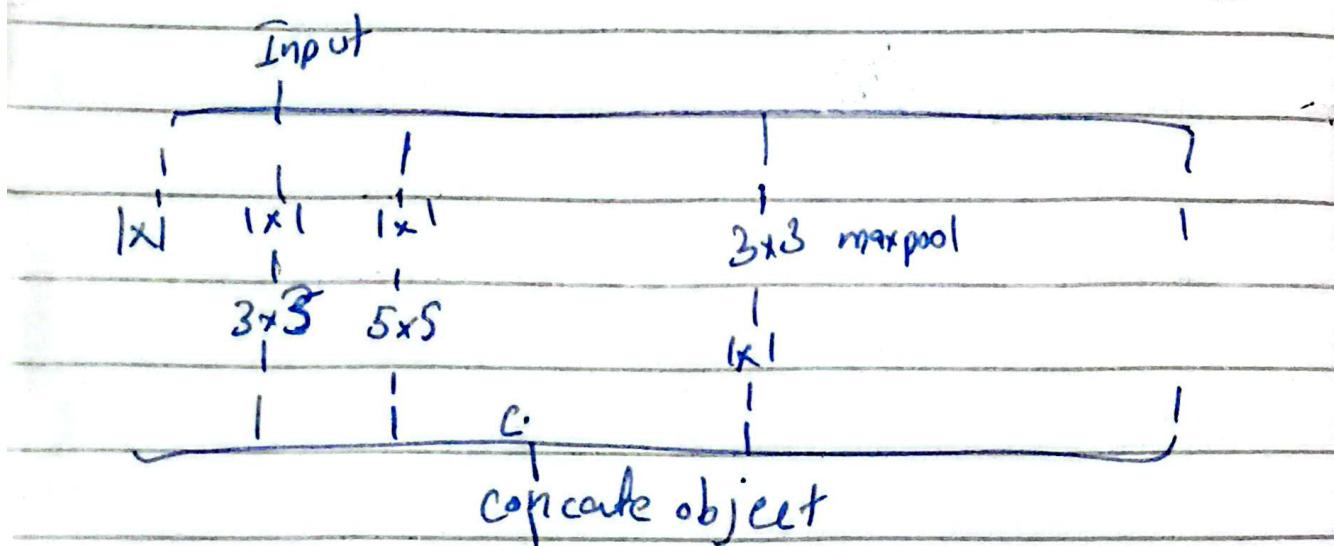
First reduce channel, then detect medium-scale patterns.

3) 1x1 conv \rightarrow 5x5 conv

Reduce channel and detect large scale.

4) 3x3 max pooling \rightarrow 1x1 conv

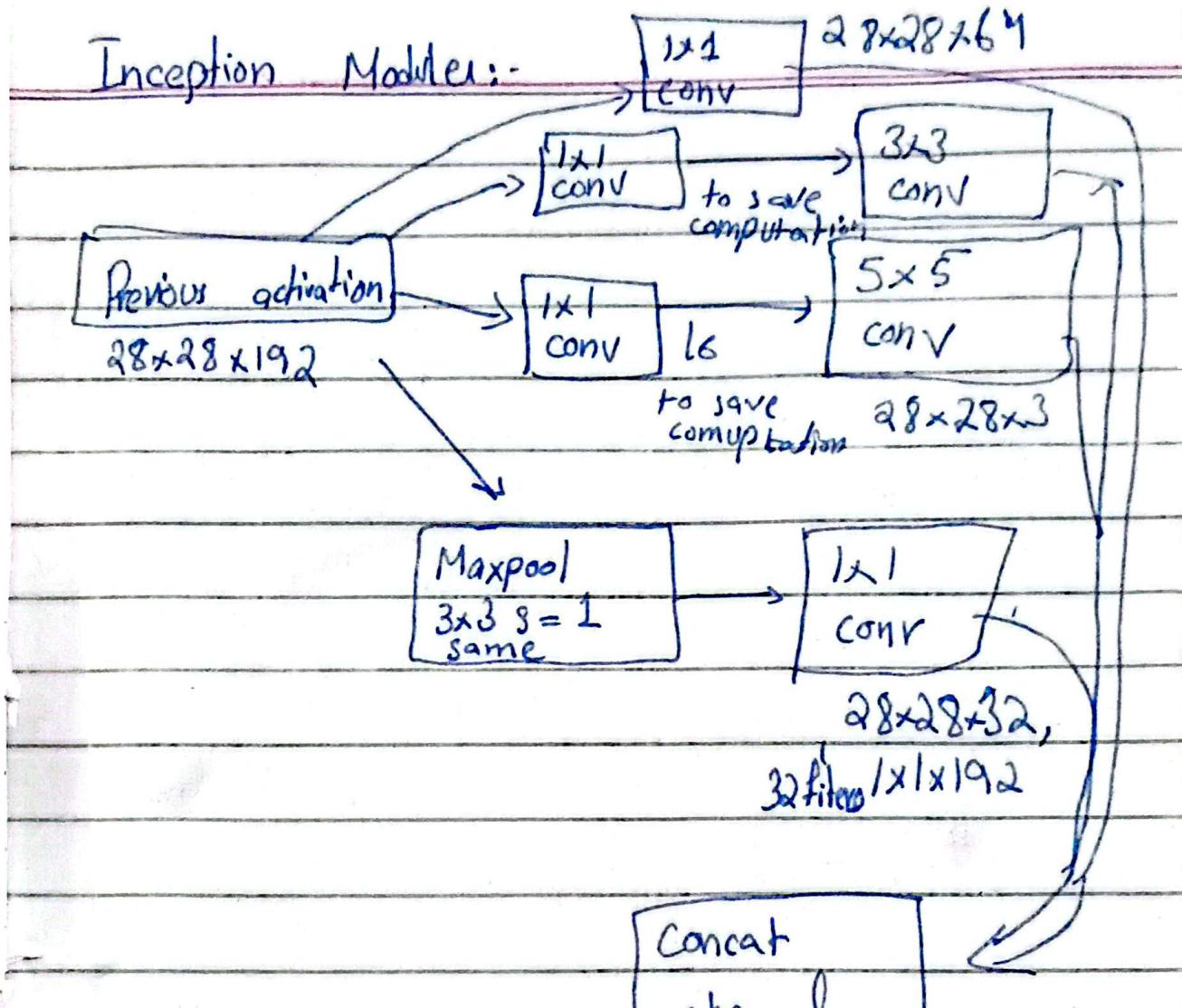
Capture background, then mix channels



When Use:-

When we can't decide use 1x1, 3x3, 5x5 convolution then we use Inception Network motivation

Inception Modules:-



$6 \times 6 \times 32$

$d = \text{filters}$

$w \times h \times \text{filters}$

$w \times h \times d$

Mobilenets:-

⇒ efficient, can use even in mobiles

Normal Convolution = computation cost 2160

$$6 \times 6 \times 3 \quad * \quad 3 \times 3 \times 3 = \begin{array}{|c|c|c|} \hline & 1 & \\ \hline 1 & & \\ \hline & & \\ \hline \end{array} \quad 4 \times 4$$

Depthwise convolution (RGB)

Red

$$6 \times 6 \times 3 \times 3 \times 3 = \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \quad 4 \times 4 \times 3$$

$$\cancel{6 \times 6 \times 3 \times 3 \times 3} \text{ computation cost} = 540$$

What happens

* this multiply single individual layers (3 RGB one by one) with a individual filter (red for red filter, blue with blue, green with green).

* This individual layer multiplication cause in reduce computation cost.

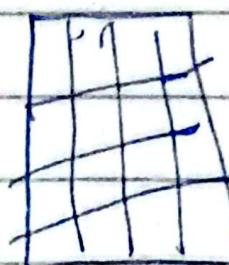
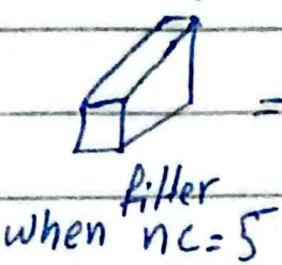
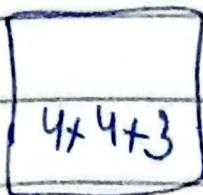
* Thus we can use this in mobiles.

* After this it again subject it to pointwise convolution.

result



further multiply to point wise



$4 \times 4 \times 1$

$\curvearrowright 3 \times 4 \times 5$

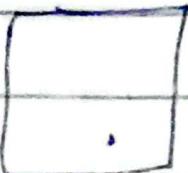
Pointwise

$$\text{Computation cost} = \text{filter param} \times \text{filter positions} \times \overset{\text{no. of}}{\text{filter}} \\ = 1 \times 1 \times 3 \quad 4 \times 4 \quad 5$$

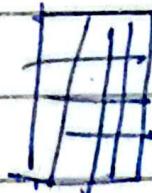
Computation = 240

Mobile Net use

Depth wise

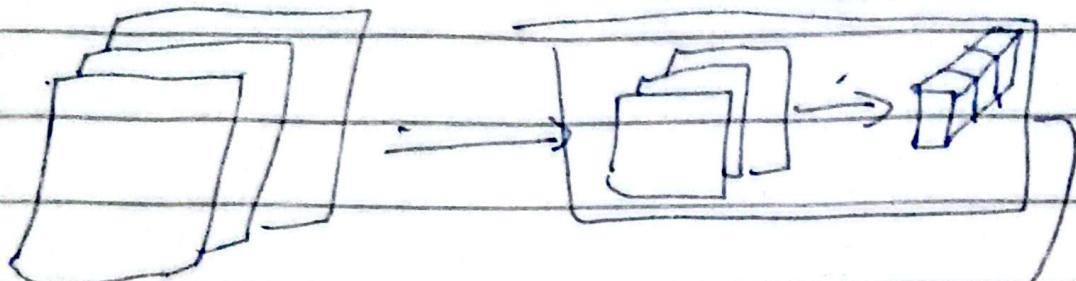


→ Convolved * pointwise =



MobileNet Architecture :-

13 times



$4 \times 4 \times 3$

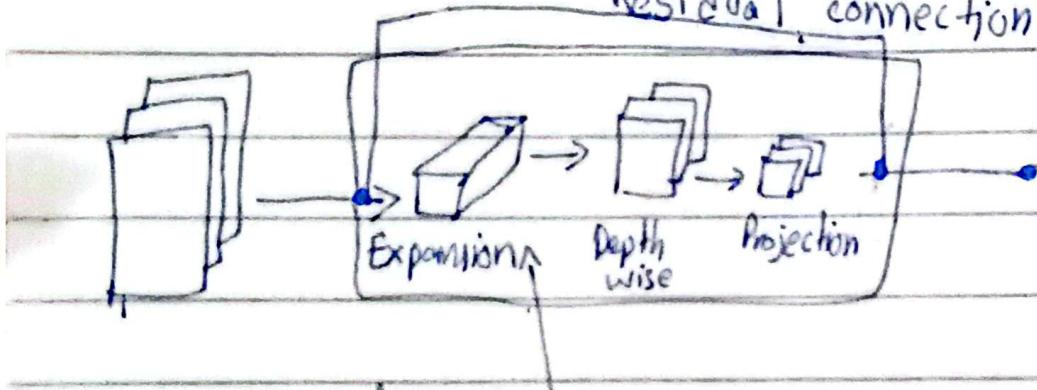
Pool \rightarrow FC \rightarrow softmax.

Improvement in it.

There are two improvements:

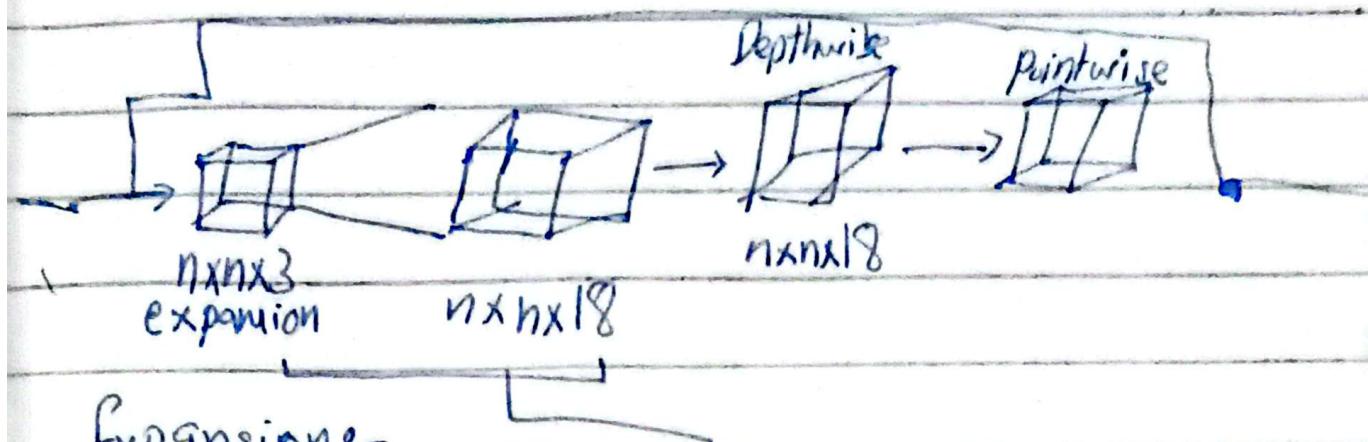
1) Addition of Residual Connection

Residual connection



2) Addition of expansion connection

Residual connection



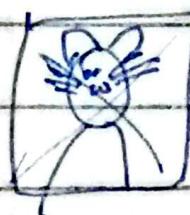
Expansion:-

increase the size

Use of depthwise and pointwise cause decrease in computation

To Efficient Net

Baseline



We can find an efficient way by efficient net to up and down the ($w \times h \times d$) for a specific device.

Transfer learning

is a technique of taking a CNN that has already been trained on a large dataset (like Imagenet with millions of images) and adapting it for a new, often smaller task.

i) Feature extraction

Freeze most of all pre-trained CNN layers

DATA
AUGMENTATION

This same as changing axis, scale size, Rotation, Mirroring is data augmentation.

Module 3 Object Detection

Object localization

localization is to detect where the car is

Landmarks

If we want to detect hand from the body we will train the machine by providing it the picture, coordinates and bounding box for its training e.g.

detect the eye

L₁...L₂, for detecting a corner

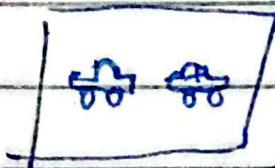
L₃...L₄, for detect another corner

⋮

Car detection example

Training set

X



[car]

1

[car]

1

[car]

1

[bg]

0

[dog]

0

Running or sliding windows require computation
but we have solution

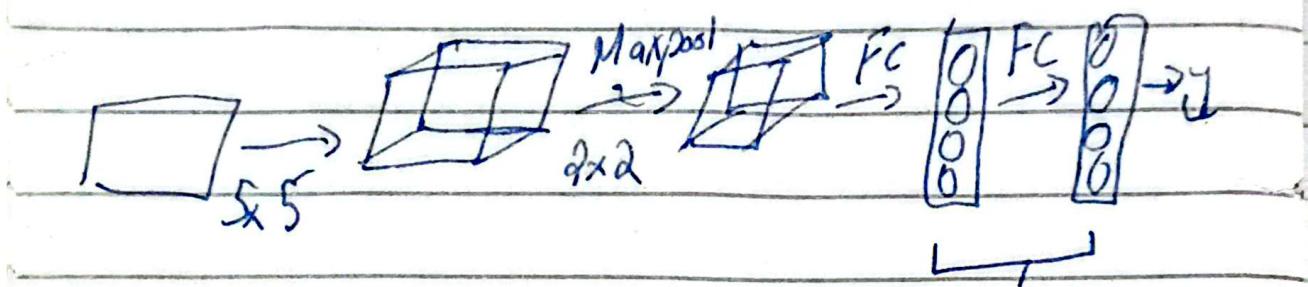
Solution:

↑ Convolution Implementation of sliding windows

$2 \times 2 \times 40$

24×24

First turning FC layers to convolution layer



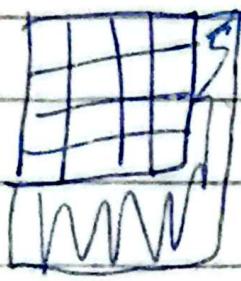
How to make them
conv layer.



IMPLEMENT of sliding windows

(FC) FC has been converted in convolution

\because ~~with~~ padding



(6x16 + 3)



12x12x16

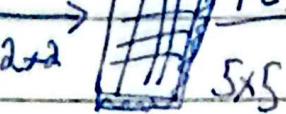
Maxpool



3x2

40

Maxpool



5x5



1x1

object

↓

FC

↓

1x1

↓

1x1

but this process is not ^{too} accurate.

It don't output correct bounding boxes.

Solution:-

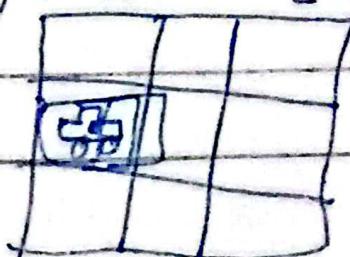
Bounding box Predictions.

\Rightarrow Yolo algorithm

\Rightarrow It paints the midpoint of target bounding box.

\Rightarrow treat detection as a single regression problem rather than running classifier on many regions.

\Rightarrow Yolo allows each grid cell to predict multiple bounding boxes for different object shapes

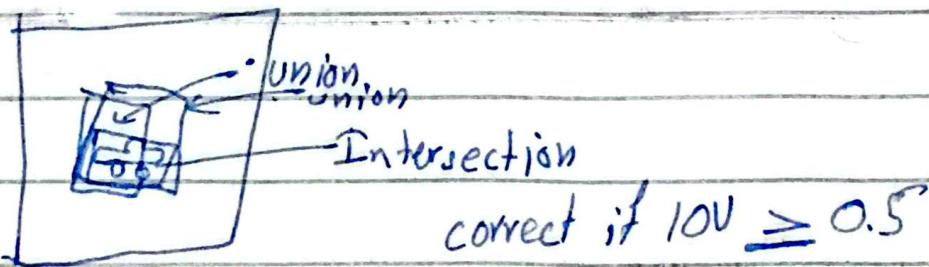


OD

How to check algorithm is working well

- Intersection over union

$$= \frac{\text{size of intersection}}{\text{size of union}} \quad \therefore \text{correct if } \text{IOU} \geq 0.5$$



It shows us how correct bounding box is, it should be 0.5, higher the number

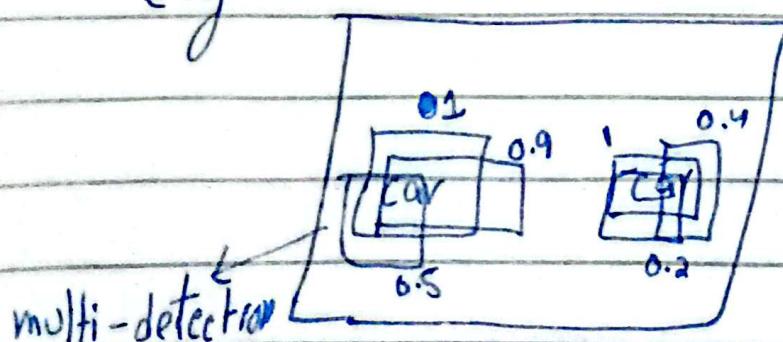
Greater accuracy is

IOU is the measure of overlap b/w two
bounding boxes.

TH Non-max Suppression:- (Help to work even better)

(It prevents multiple detection of same object:-

e.g



Non-max suppression prevent this.

(as 1)

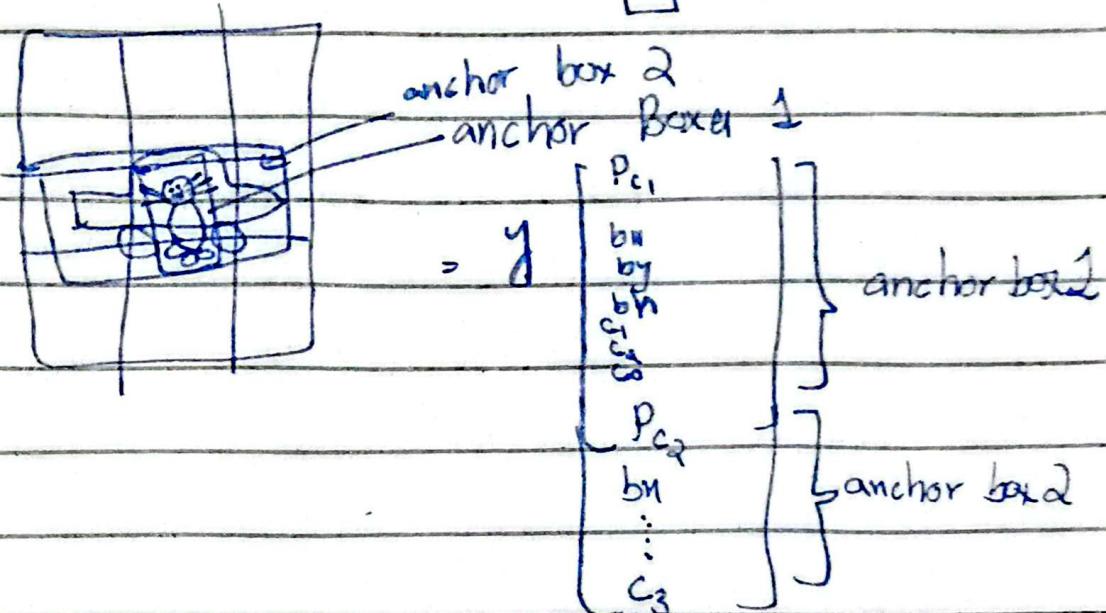
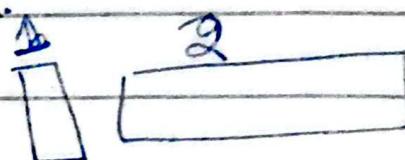
This cause to take max probability and other boxes are suppressed with less probability.

To Anchor Boxes:-

Object detection

One of problem is that one grid detect only one box.

e.g.



Previous with one box or object

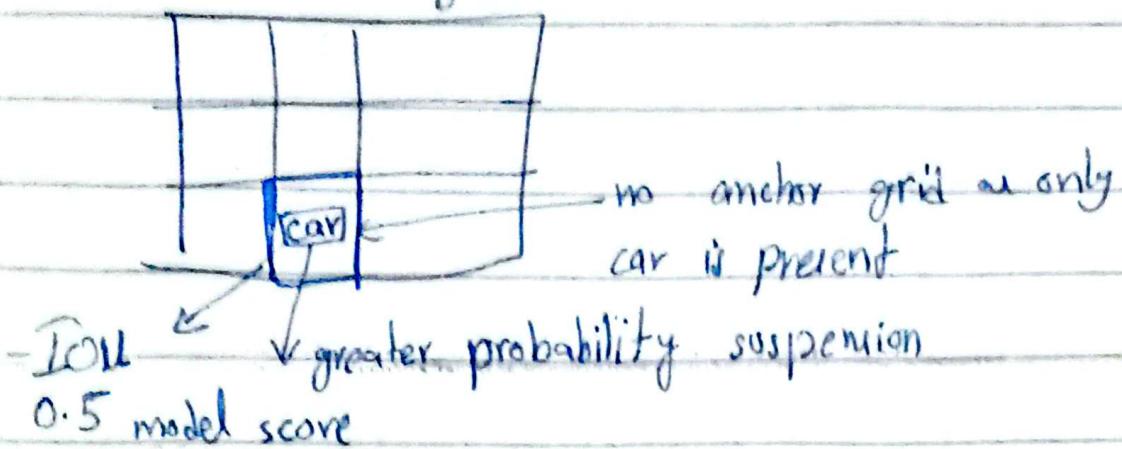
$3 \times 3 \times 8$

now \longrightarrow 2 anchor boxes.

$3 \times 3 \times 2 \times 8$

Put all the things (IOU, Anchor, Non-mansuped) in Yolo

Yolo algorithm



Region proposal cell (often less used)

This allows the sliding window where it finds the target rather than the sliding over the whole image.

This algorithm is slow

FAster R-CNN:

Use convolution network to propose but slower than Yolo algorithm

Semantic and Segmentation with U-Net

Semantic segmentation

is a computer vision task (car-driving) where the goal is to assign a class label to ~~class~~ every single pixel in an image.

Search object vs semantic segmentation

Useful

Chest X-Ray, Brain MRI

~~It this~~

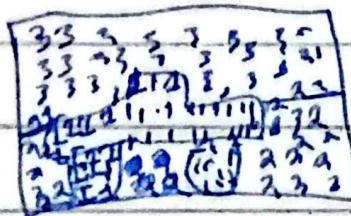
Working :

It assigned the pixel of specific thing (like car, road, building) with number

1 for car

2 = road

3 = Building



and also give color to ever pixel

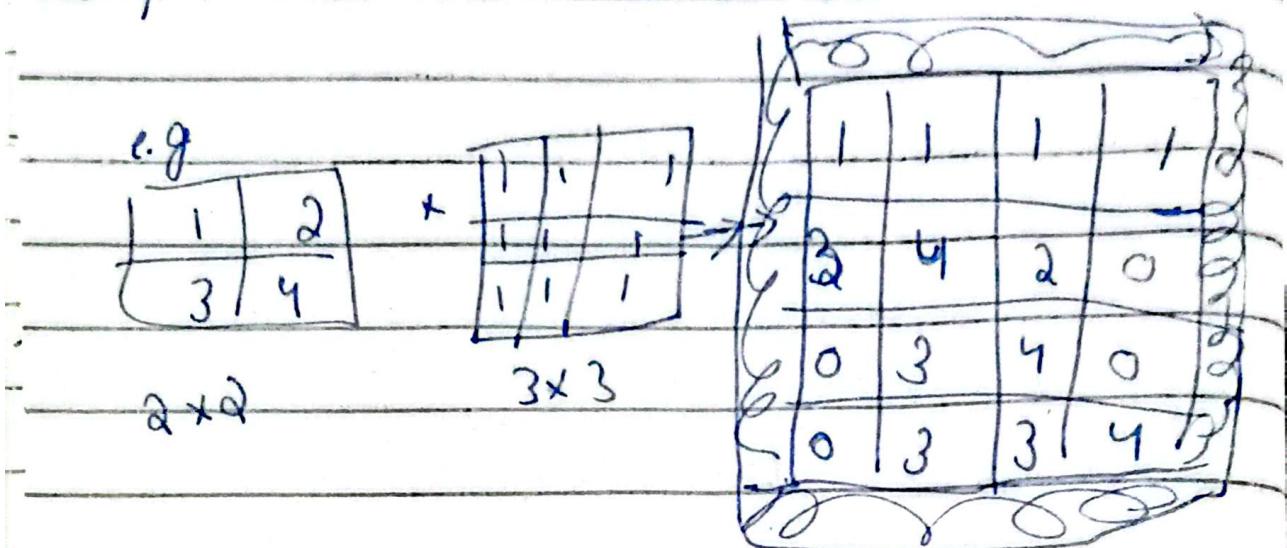
red = 1

Blue = 3 = Building

2 = Black = road.

Transpose convolution:

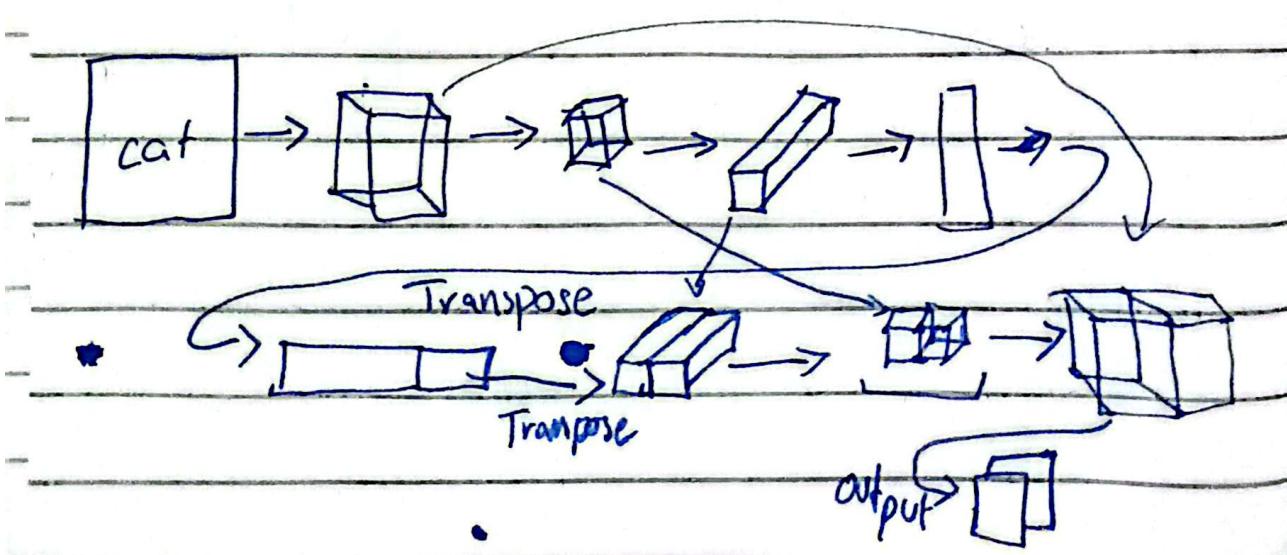
is a way to generate large output from the smaller one



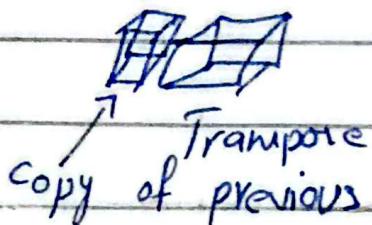
padding = $P=1$, stride $s=2$

When two box overlap we ~~multiply~~ add then

U-Net Architecture intuition::



After taking transpose it also merge the previous copy in UNet



Week #4

Face Recognition - VS face ^{verification} ~~recognition~~

Verification:-

- Input Image, name/ID.
- Output whether the input image is that of the person claimed person (Detect + comparison)

Recognition.

number of

- Has a database of K persons
- Get Input
- Output ID if the image is any of K persons.

[Image + compare from db + get ID]

One shot learning (algorithm recognize in one view)

* Like ability of algorithm to recognize by any image by getting one image of that person.

* Recognize different axis of face angle by trained on only one image of that face

$$\text{If } (\text{image 1}, \text{img 2}) \leq \tau \\ > \tau$$

Siamse network.

is a type of neural network architecture that contains two or more identical subnetworks that share the same weights and parameters.

a) Each input is passed through the same neural network - This network extracts feature vectors from both inputs.

b) The network calculates a distance or similarity score between the two vectors.

c) If the inputs are different, the distance

should be large.

Triplet loss

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

Def:-

Triplet loss function is the loss function mainly used in metric learning and tasks like face recognition, image retrieval, or signature verification.

It works:-

It uses three inputs

- 1) Anchor
- 2) Positive (example similar to anchor)
- 3) Negative (example ~~similar~~ not similar to anchor)

Mathematical formula:-

$$\text{Loss} = \max(0, \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha)$$

Example:-

Suppose: $d(A, P) = 0.5$, $d(A, N) = 0.9$

$$\alpha = 0.3$$

$$L = \max(0, 0.5^2 - 0.9^2 + 0.3) = \max(0, -0.1) = 0$$

Choose triplets that are hard to train.

T# Face verification and binary classification.

$$\hat{y} = \alpha \left(\sum_{k=1}^{128} |f(x^{(i)})_k - f(n^{(j)})_k| + b \right)$$

$$\text{chi square} = \frac{(f(x^{(i)})_k - f(n^{(j)})_k)^2}{f(x^{(i)})_k + f(n^{(j)})_k}$$

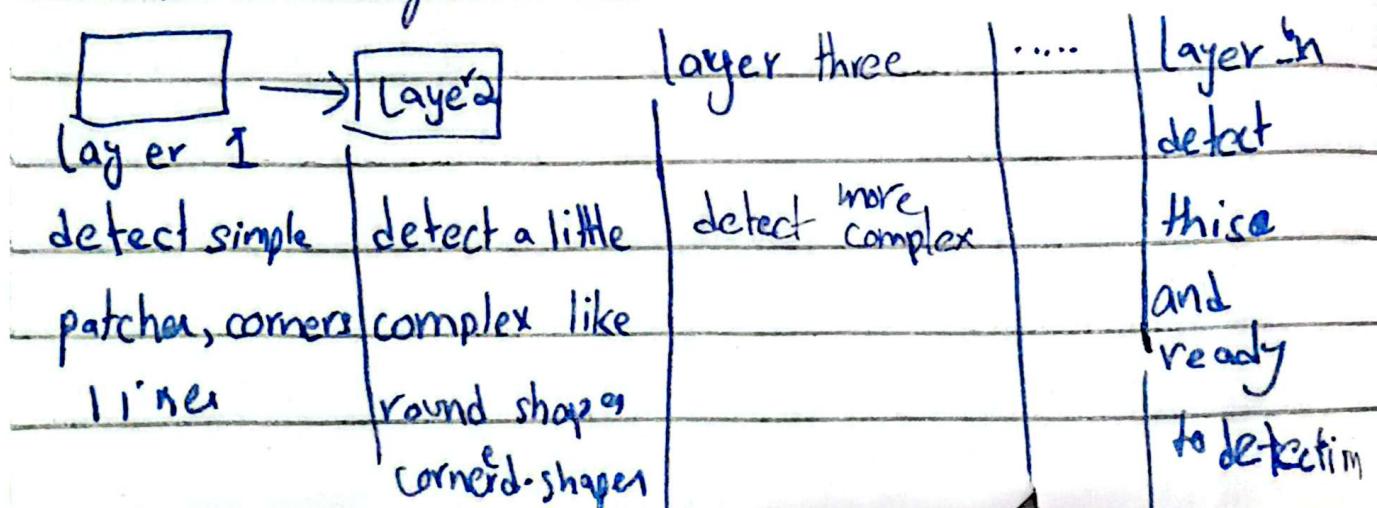
Binary classification

is the more general ml task where the system classifies input data into one of two categories or classes.

Input: A single data sample (Two faces)

Output: One of two classes { Yes or No }

T# Neural style transfer



Deep Convolutional Neural Networks.

These use convolution layers to automatically and adaptively learn spatial hierarchies of features from simple edges in early layers to complex objects in deeper layers.

T# Neural style transfer cost function - of Generated image

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

2) Use gradient decent to minimize

$$G = G - \frac{\delta J(s)}{\delta G}$$

T# How to define content cost function

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

Say you use hidden layer L to compute content cost.

→ Use pretrained ConvNet

→ Let $a^{[L][C]}$ and $a^{[L][G]}$ be the activations of layer L on the images.

If $a^{(L)(c)}$ and $a^{(L)(g)}$ are similar then they have similar content.

Hence, this is used to detect different content

$$J(G, G) = \frac{1}{2} \|a^{(L)(c)} - a^{(L)(g)}\|^2$$

$(G_s - G_G)^2$
instead of difference :

$$(G_s - G_G)$$

The style cost function should be:

$$J_{\text{Style}}^{(L)}(S, G) = \frac{1}{n_H} \sum_k (G_{kk}^{(L)(S)} - G_{kk}^{(L)(G)})^2$$

$$= \frac{1}{(2n_H^L n_W^L n_C^L)^2} \sum_k \sum_k (G_{kk}^{(L)(S)} - G_{kk}^{(L)(G)})^2$$

If $a^{[L](c)}$ and $a^{[L](g)}$ are similar than they have similar content.

Hence, this is use to detect different content

$$J(c, g) = \frac{1}{2} \|a^{[L](c)} - a^{[L](g)}\|^2$$

$(G_s - G_g)^2$
instead of difference:

$$(G_s - G_g)$$

The style cost function should be:

$$\underset{\text{Style}}{J^{[L]}(s, g)} = \frac{1}{n_H n_W} \sum_{k=1}^K \|G_{kk}^{(s)} - G_{kk}^{(g)}\|^2$$

$$= \frac{1}{(2n_H n_W n_C)^2} \sum_k \sum_{k'} \left(G_{kk'}^{(s)} - G_{kk'}^{(g)} \right)^2$$

Style cost function

Style:

as correlation between activations across channels. (corr) between style of patches

Style matrix

Let $a_{ijk}^{(l)} = \text{activation at } (i, j, k)$. $G^{(l)} \in \mathbb{R}^{H \times W \times C}$

Note: Neural style transfer uses this style matrix to generate new images that combine the content of one image with the style of another.

Style Matrix: Encodes style features for style transfer.

$$G_{kk'}^{(l)(s)} = \sum_{i=1}^n \sum_{j=1}^n a_{ijk}^{(l)(s)} a_{ijk'}^{(l)(s)}$$

$$G_{kk'}^{(l)(s)} = \sum_{i=1}^n \sum_{j=1}^n a_{ijk}^{(l)(s)} a_{ijk}^{(l)(s)}$$

$$J_{\text{style}}^{(l)}(s, G) = \| G^{(l)(s)} - G_l^{(l)(s)} \|_F^2$$

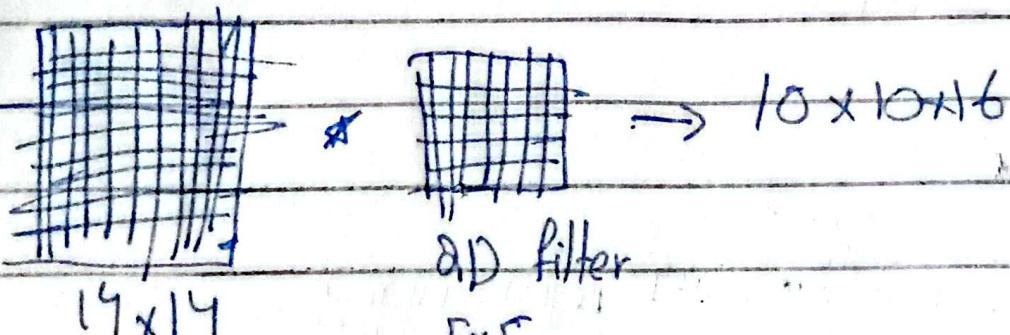
$$= \frac{1}{(2n_H^{(C)} n_w^{(C)} n_c^{(C)})^2} \sum_k \sum_{k'} (G_{kk'}^{(C)(S)} - G_{kk'}^{(C)(G)})^2$$

Overall style cost function

$$J_{\text{style}}(S, G) = \sum_l \lambda^{[l]} J_{\text{style}}^{[l]}(S, G)$$

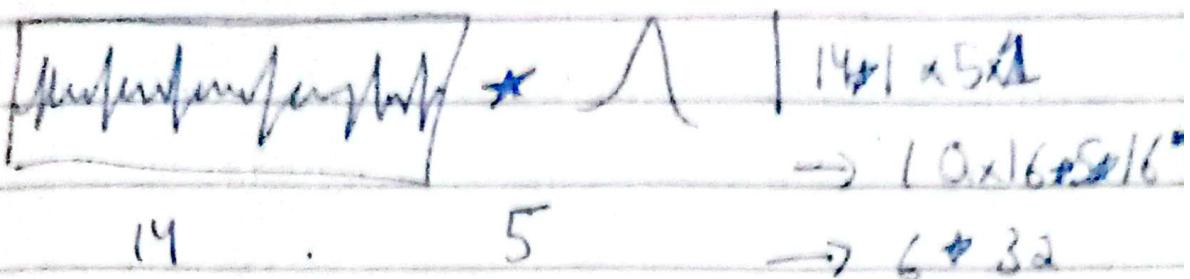
$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

1D, 2D and 3D Images.



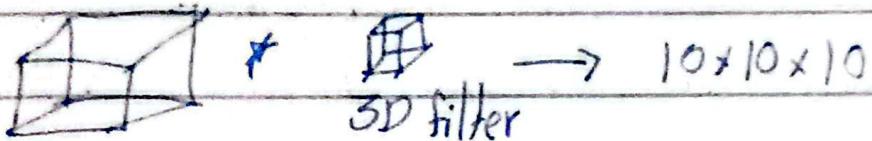
2D

When we use ecg (1D) Data



What about 3D data

CT scan, organs scan



Sequence model:-

Sequence model is one which trained on sequential Model.

* DNA sequence

* Speech Recognition

* Music Generation.