

04. Feb. 2018

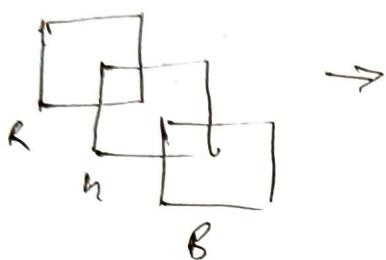
> Vanilla architecture - there is no regularization

(bog) Mnist $\rightarrow 28 \times 28 \times 1$ channel
Label
(Rip) CIFAR $\rightarrow 32 \times 32 \times 3$ channel

Vanilla Architecture

Basic CNN

Conv. layer



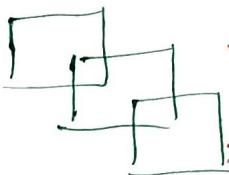
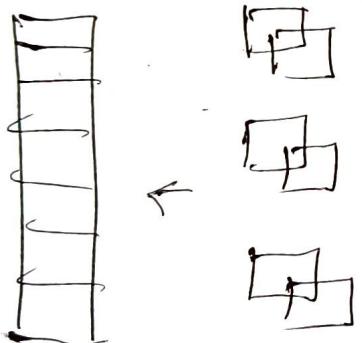
ReLU



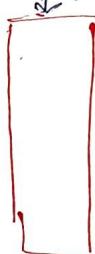
Feature Map Volume

Flatten layer

Tensor

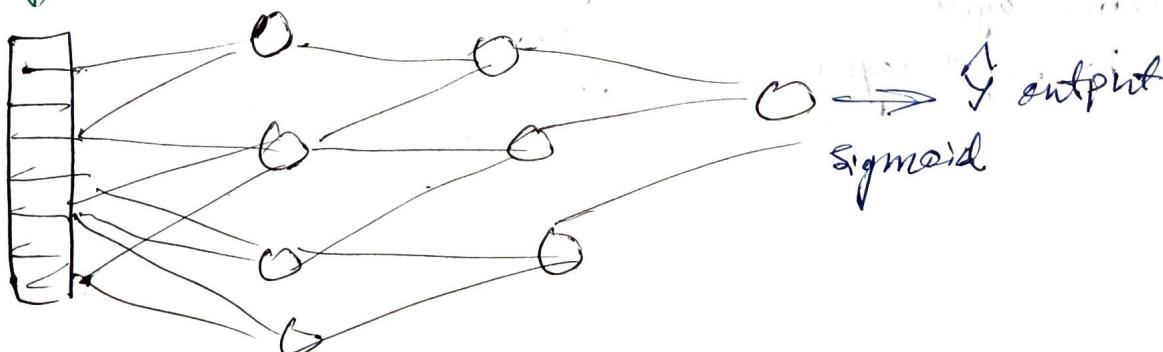


Pooling layer



If brings
data into
1-D array

ANN



ImageNet Competition

→ Researchers have used this dataset to build different CNN architectures.

- > People changes the architecture by
 - ① convolution layer
 - ② filter size
 - ③ stride
 - ④ Padding
 - ⑤ Fully connected layers (ANN)
 - ⑥ activation function
 - ⑦ Dropout
 - ⑧ Batch Normalization

History of CNN

- ① LeNet - 5
- ② Alexnet
- ③ VGG (VGG16 or VGG19)
- ④ Inception (GoogLeNet)
- ⑤ ResNet (Microsoft)

Note: there are many CNN architecture available.

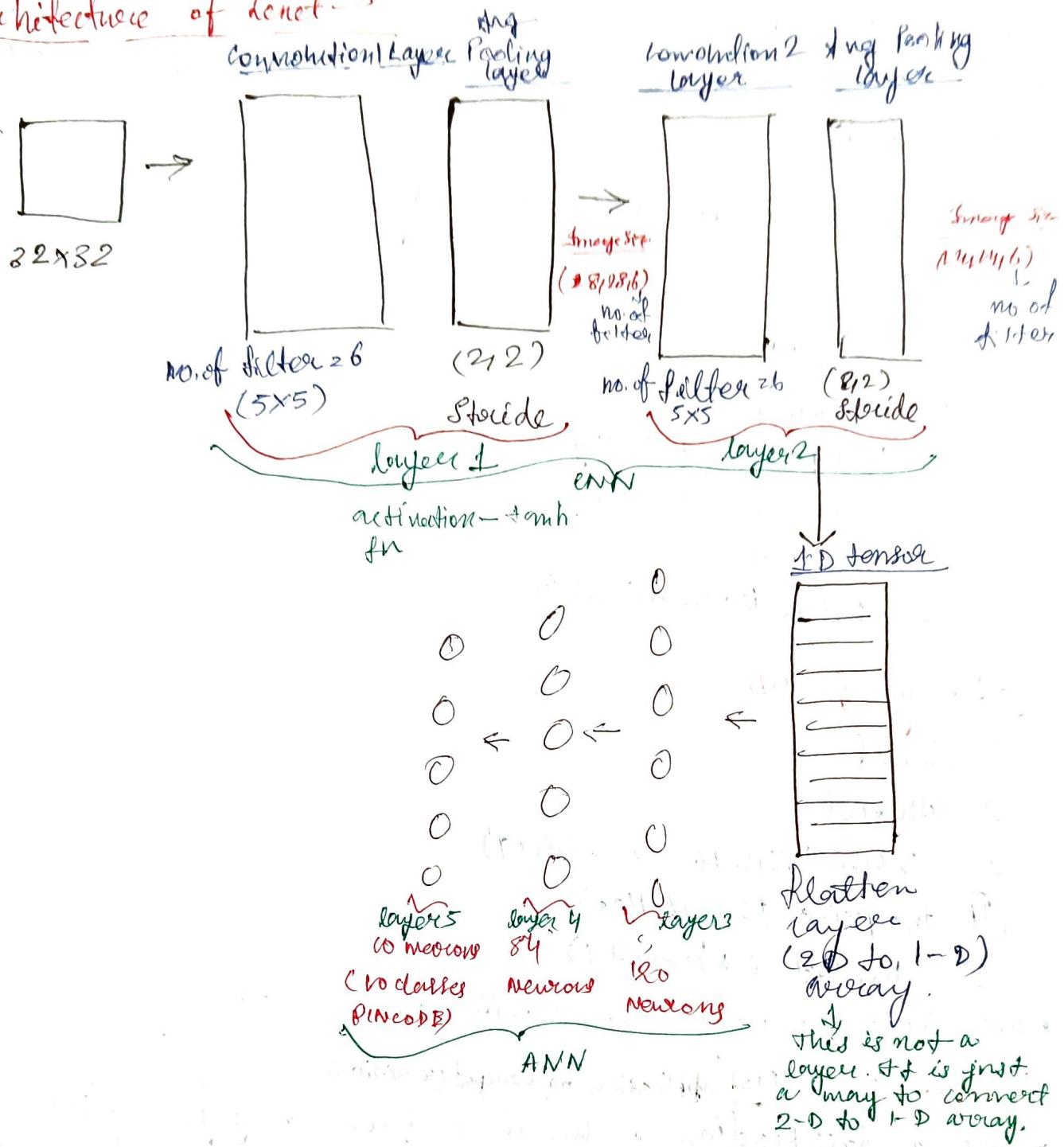
① LeNet - 5 (1st application of computer vision)

- > It's a classification architecture from a paper called
- > A gradient based learning for document classification.
- > Handwritten digit classification / character VS Navy
- > Postal Service - developed this architecture to recognize pincode.

1245 = Pincode

- > Published in proceeding of the IEB (1998)
Input size = 32x32x1 (1st nice work for this input only)
- > This paper was published by Yann LeCun, Leon Bottou and Yoshua Bengio and Patrick Haffner.

Architecture of Lenet-5



1st conv layer

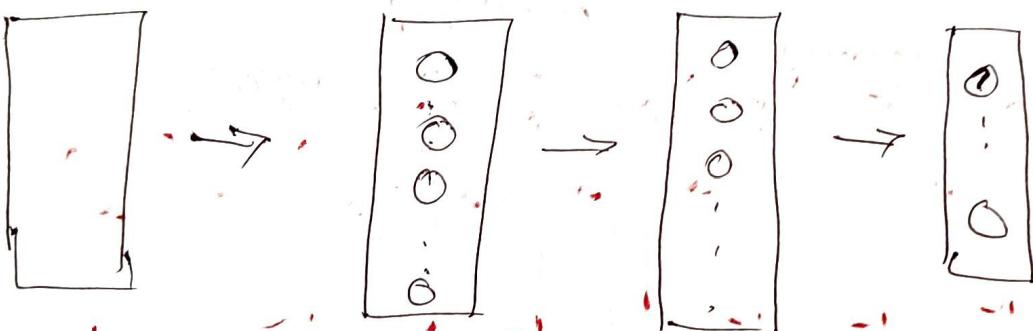
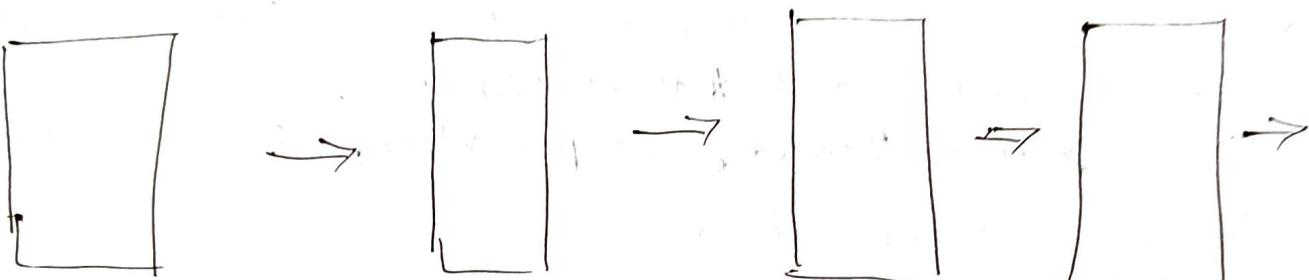
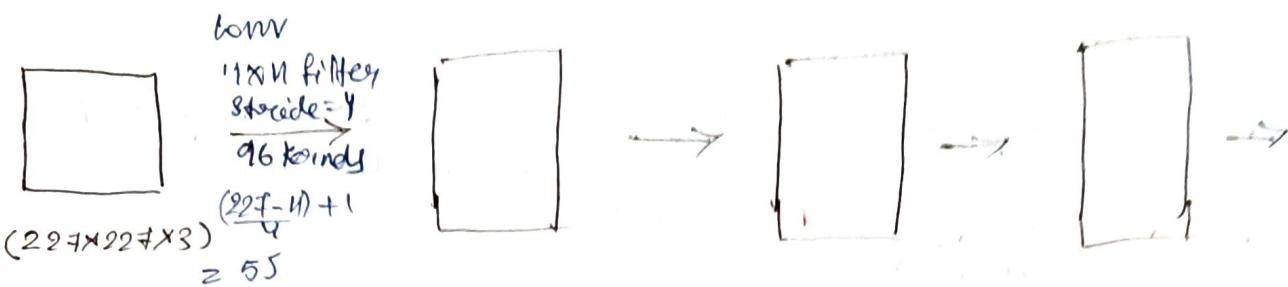
$$= \frac{(32 + 2 \times 0 - 5) + 1}{1} \\ = 28$$

$$\text{Image dimension} = (28 \times 28 \times 3)^3$$

> **RNN Performance on RNN image model be very poor because no. of layers is very less.**

② AlexNet

2(L-54)



05-Feb-2022

⑧ AlexNet (16 layers)

> It's a classification architecture / model, from Neel Ganguly. Convolutional networks for large-scale image recognition.

Application of object computer vision

① Image classification

② Object Detection

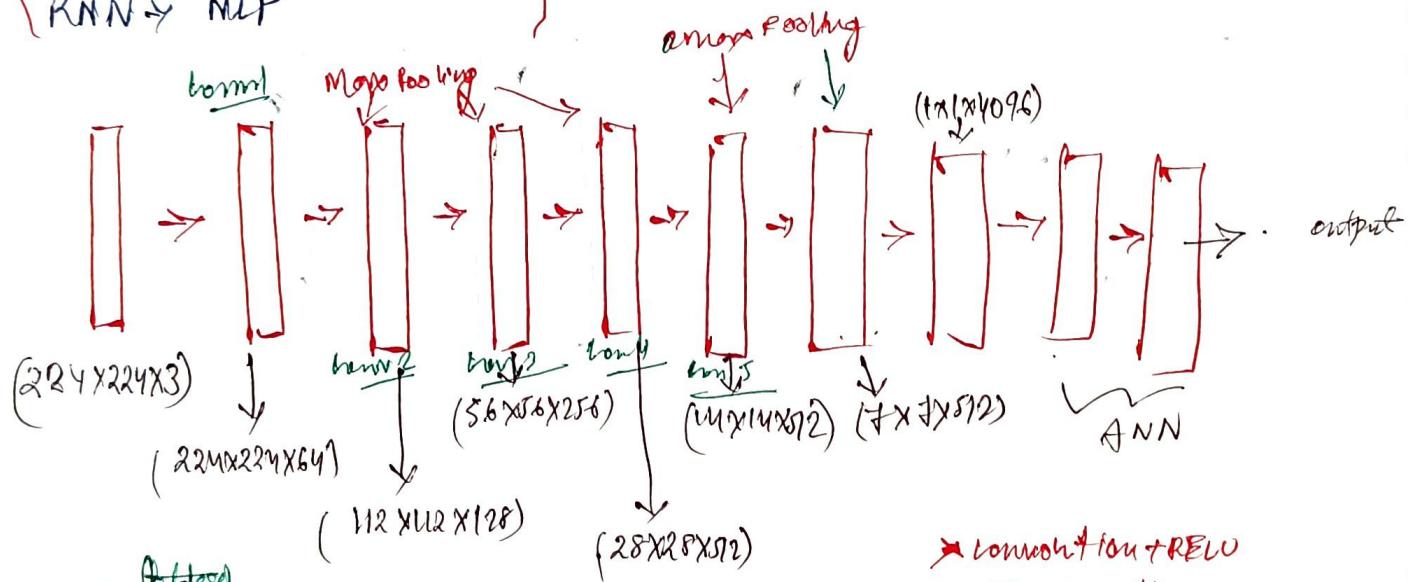
③ Face Recognition

> It is published as a conference paper at ICLR in 2015, by Karen Simonyan and Andrew Zisserman.
Visual Geometry Group - VGG → from OXFORD.

> Video is a sequence of images.

{ CNN → Computer Vision } use cases

{ RNN → NLP }



→ No. of ~~filter~~ ^{feature} is getting increased.
no. of filter

- * Convolution + ReLU
- * Max pooling
- * Fully connected + ReLU

> All convolution layer has filter size = 3×3 . But in Alexnet filter size is different for different layer, in very all stride is 1, and padding is same.

> In max pooling layer, pooling size is 2×2 and stride is 2×2 .

VGG-19

- It has 19 layers.
- So frontal for all the architecture.

Q why we use pre-trained model?

- Biggest challenge in DL is availability of data. DL and techniques are data hungry.
- Always labeled the data
- Model training time.
- To overcome the above ~~time~~ ^{time}, we use pre-trained model.

ImageNet Dataset

- gives a visual representation of images.

why?

In 2006: ^{main focus} Model and algo building to get better accuracy.

- Garbage in, Garbage out.

Bad data bad model Bad Result

Good data bad model Good Result

- After this one researcher thought to build big good database. He collected around 1.4M images, and labelled the images and their features like colors, etc.
- He used crowd sourcing + scratch sourcing to collect the data. Using captcha technique
- Collected data is known as ImageNet. It is the largest dataset of the image classification in the world as of now.

ILSVRC / ImageNet challenge

→ ImageNet Large Scale Recognition challenge (ILSVRC)

1.4 M data - 20,000 class

1 million → ImageNet competition
1000 class

Machine Learning

Fannodysis $\xleftarrow{\text{SFT}} \text{ML}$

2010 → 28% error rate } ML model
2011 → 25% error rate

Finally 2012 → Reduce the error to 16% by DL

Jeffrey Hinton, he propose deep learning based architecture called AlexNet and trained on GPU

In 2013 → Error rate was 11.7%.

In 2014 → VGG → 7.3% Error rate

Inception In 2015 → Google Net → 6.7% Error rate

In Fe 2016 → ResNet → 3.5% error rate.

Deep learning

Image classification Why to build from scratch?

> we will use pre-trained model.

eg To build a car, we use same technique to build wheel.
Because invention has been done already.

Transfer Learning

> use CNN part and drop ANN part.

→ feature extraction

> If we train convolutional layers, then it is called fine tuned.
When dataset is completely new.

11-Feb-2022

Computer Vision

Filter / Kernel Feature Extractions

- Use is to extract features

Eg a Bird & DFF How to decide whether it is name or not using filters?

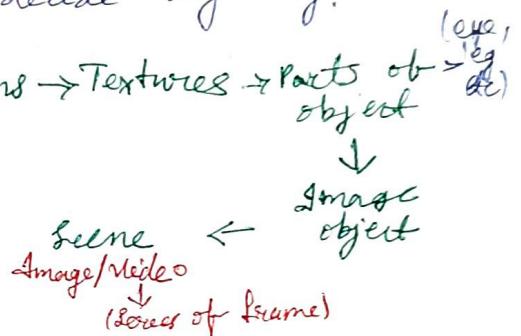
Pixels

> Smallest component

> After grouping of pixel, we can infer something.

> Only based on pixel, we cannot decide anything.

Pixels → Edges / Gradients → Patterns → Textures → Parts of object



> 5 layers to get scene

Image RGB $\geq h \times w \times 3$

Purpose of filter

{ Small to Big
Pixel to Scene of object }

Big to Small

Image → object → Part of object → Textures → Patterns → Edges → Pixels

Dimensional

This can be anything

200 × 200 × 3
RGB

200 × 200 × 12

12 different
channel / feature map
Group of channels

Image
colours

Feature maps

No colours

- Decision will be made on shape
not on colour.

12 filters \rightarrow input

$200 \times 200 \times 12 \rightarrow$ output

why 12 filters?

1st decide Kernel then no. of kernel

$n \times n$

which is best?

4×4

Why 3×3 best?

3×3

\rightarrow has less parameters.

5×5

7×7

Ex

① 36×36 image \leftarrow (input)

$3 \times 3 \times 12$
(filter) 12 times

$P = 0$ filter

$$(n - 2p + f) + 1 = (36 - 0 - 3) + 1 = 34$$

~~$p = f + 1$~~
 $\text{output} = 34 \times 34 \times 12$

$3 \times 3 = 9$ (2)

$$\text{Params} = 9 + 9 = 18$$

less important information
with 2 3×3 filters
compare to single 5×5
filter.

\rightarrow If I use 3×3 , 2 pixel will be lost.

② input $\Rightarrow 36 \times 36 \times ?$

Filters $\Rightarrow 5 \times 5 \times 12 \Rightarrow$ 12 channel

$$(n - 2p + f) + 1 = (36 - 0 - 5 + 1) = 32$$

$$\text{output} = 32 \times 32 \times 12$$

12 \rightarrow no. of filters, It is a hyperparameter.

$5 \times 5 = 25$ (single)

$$\text{Params} = 25.$$

more information
lost with single
filter

$$360 \times 360 \times 3$$

single convolution operation

$$360 \times 360 \times 2$$

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

$$5 \times 5 \times 1 \quad (225) \text{. Parameters}$$

$$358 \times 358 \times 1$$

$$356 \times 356 \times 1$$

$$358 \times 358 \times 2$$

Local convolution operation

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

$$356 \times 356 \times 1$$

$$\text{Params } 2 \times 2 \times 3 \times 3 \times 1 = 18$$

> so, 3×3 is better than 5×5 based on less information loss and parameters.

Information loss

> when we are using 2 filters, information lost can be stored in feature map. But in 5×5 we directly loss information.

Padding

> Trunc

Padding

> added up values

> Black pixels

> Adding padding so that we can do more convolutional.

Ex Input $h \times w \times 3$

Output $1 \times 1 \times 256$ (Before flatten) (we need in this shape)

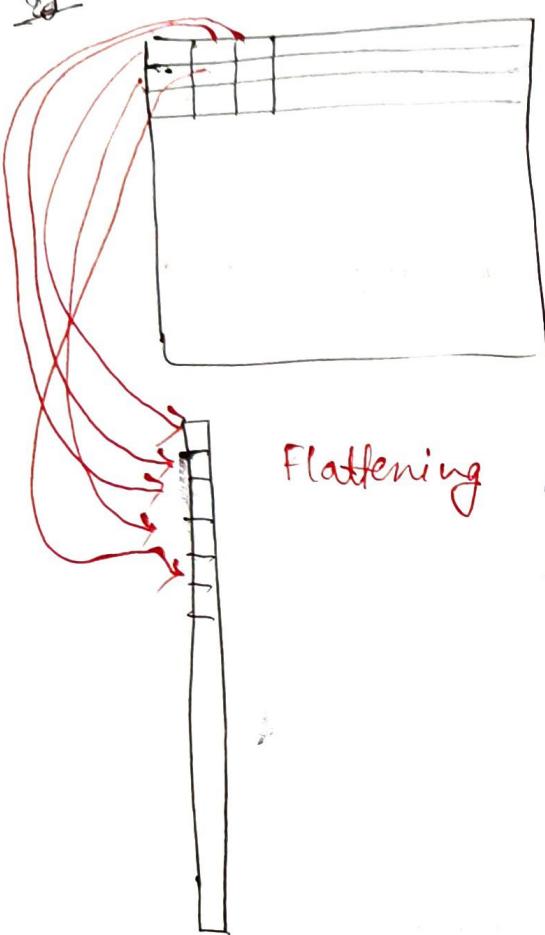


> convolutional means moving of the kernel on top of the image to extract features.

ANN

→ It learns by learning more complex features.

→ 3. S → Distinguish identification can't be done from computer image.



* Appear all the friends medically

Flattening

- > If we directly flatten them it will not hold shape. It is of no use.
 - > If we have to perform convolutional operation to hold the shape.

Can we make flatten in DPL?

~~the~~ sense gentle flatten sense.

Yes it can be.

> using too much dense is bad,

Sente

One
decapent

Sense

Hoffman

~~genetic~~

flattened
seems (no) sprayer.
soffner

softmax

CNN Architecture

- ① Lenet
- ② Alexnet
- ③ VGG

} All has Fully connected layers.

Q Can we do without fully connected layer?

convolutional

Input

Conv

Conv Padding

Conv

Conv

Conv

Max Pooling

Flatten

FC - Fully connected

Softmax

} Structuring 2,5,6,10,2,5,6,10 (Image size)

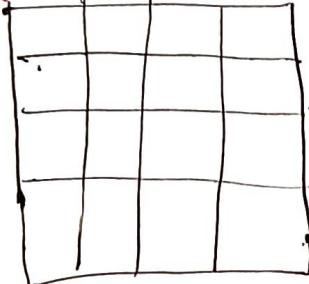
$1 \times 1 \times 10$ (classing size)

} classing

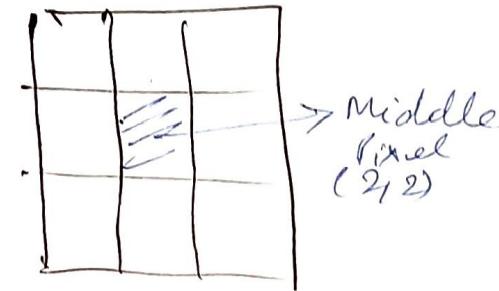
- We can have only one flatten.
- Flatten disrupts the initial information.

Q Why odd filter is better?

Eg Any



→ There is no middle



Total Pixel = 9

Total Pixel = 16

→ Symmetry issue

↓
Due to symmetry we use odd size filter.

Eg $36 \times 36 \times 3$ (Image size)

8×8 (filter size)

29×29 (odd pdf)

- Center has is getting focused more.
- Hardware wise also, 3×3 it's faster and good.

> all our assumptions in ex is that object is always in the middle.

> so prediction - transformation is better and more robust with batch as better.

⑧

> we use max pooling to avoid lots of layers.

> More the no of neurons, more will be the complexity