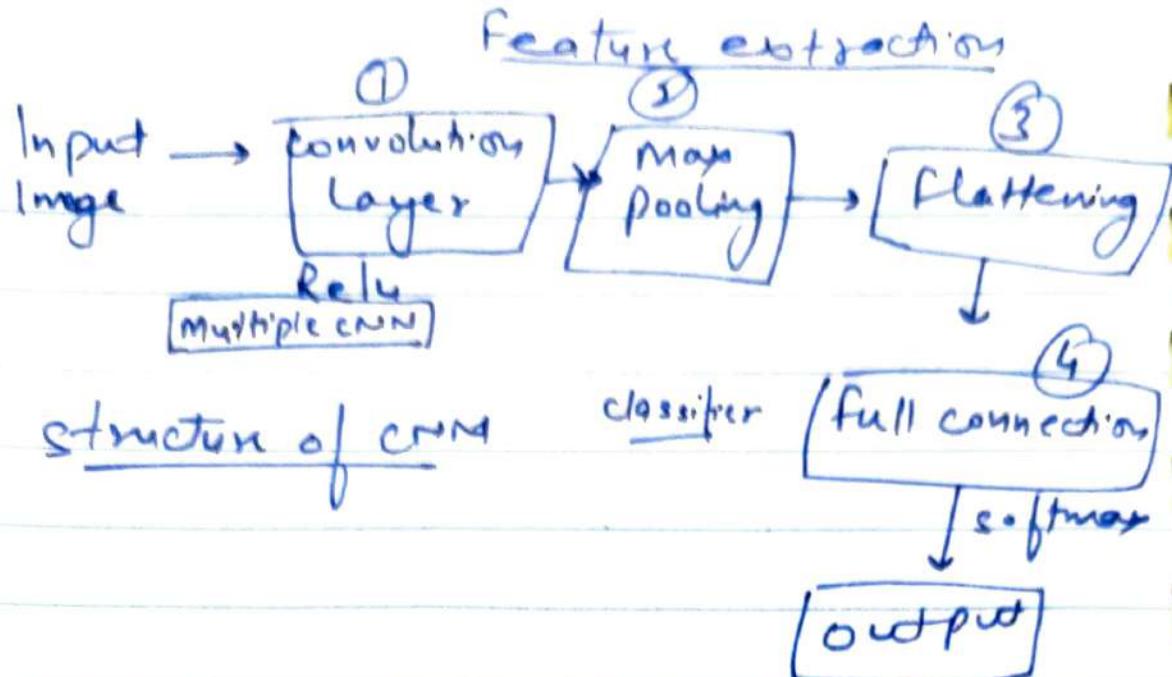


- Convolutional neural network
- CNN are widely used in image recognition
 - in image classification, objects detection, recognition face etc.
- CNN image classification takes an input image, process it and classify it under certain categories
- CNN is another type of neural network that can be used to enable machine to visualize things and perform tasks such as image classification, image recognition, object detection etc
- Image classification is the task of taking an input image and outputting a class (a cat, dog etc) or a probability of classes that best describe the image

CNN is a specialized type of neural network model designed for working with image data



structure of CNN

classifier

Input Image → CNN → Output label

(Image class)

Output number that

describe the

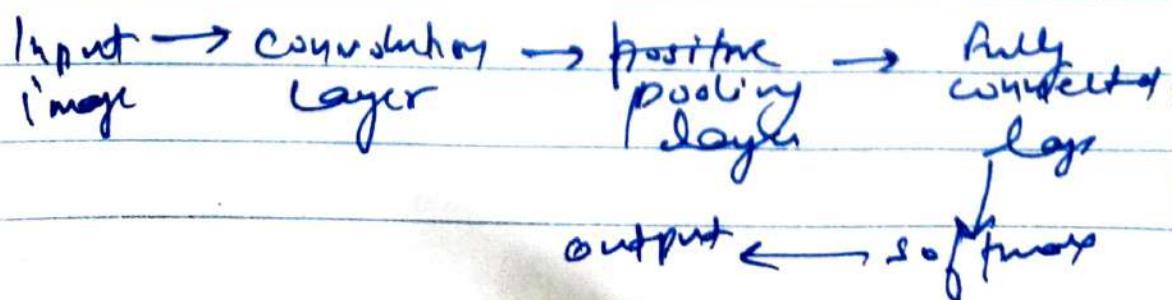
probability of the

image being a
certain class

* Array of number

In a 3D image every pixel will have 3 values separate channel for red green & blue

Each 0 - 255 value describes the pixel intensity at that point



we want the computer to do is to differentiate between all the images it's given.

- for that output computer perform image classification by looking for low-level features such as edge curves and then build up to more abstract concept through a series of convolutional layers
- In CNN the input image pass through a series of convolutional layers, pooling (downsampling) layers and full connected layers and finally produced the output which can be a single class or a probability of classes that best describes the image

Feed forward neural network can learn a single feature representation of the image but in case of complex image, NN will fail to give better predictions, this is because it can't learn pixel dependencies present in the images

CNN can learn multiple layers of

of feature representations of an image by applying different filter / transformation.

In CNN, number of parameters for the network learn is significantly lower than multilayer network.

Convolution operation: It's a linear application of a smaller filter to a larger input that result in an output feature map.

~~steps~~ Convolution layer - This layer perform an operation called a convolution hence the network is called convolution neural network

- It extract features from the input image
Convolution is a linear operation that involves the multiplication of a set of weights with the input
- The technique was designed for 2D input. The multiplication is performed between an array of input data and a 2D array of weight called filter / kernel

- Size of filter is smaller than the input data
- multiplications performed between a filter size patch of the input data and the filter is a dot product, which is then summed and generate a single value
- The filter is applied systematically to each overlapping part of filter size
 - d patch of the input data, left to right step to bottom
- If the filter is designed to detect a specific type of feature in the input (image) then it discover that feature anywhere in the image
- The output from multiplying the filter with
- If 3 filter have been used then feature map have depth of 3
- The more number of filter the more accurate result

stride Stride is the number of pixels by which we slide our filter matrix

over the input matrix

stride1 :- move the filter one pixel at a time

stride2: _____ two _____

Larger stride will produce smaller feature maps

7×7

input image

3×3

filter/
kernel

5×5

output

Input image :- array of pixel values
filter set of weights

feature map :- filter is systematically applied to the input data to create a feature map

- Multiple small feature can be designed to detect multiple feature in the image.
- Value of the filter or weights to be learned during the training of the network so that most useful features are extracted from the input image to classify it
- CNN learn multiple filter in parallel it learn from 32 to 512 filter in

parallel for a given input

- Convolution layer not only applied to input data (real pixel value) but they can also be applied to the output of other layers
- stacking of convolutional layers allows a hierarchical decomposition of the input

$$n - f + 1$$

weight sharing : In CNN weights are being used to compute the output

In ANN each neuron present in the hidden layers will have separate weights for itself

Padding : In convolution operation, reduces the size of image (spatial dimensions decrease faster)

- Zero padding allows us to control the size of the feature maps
- Padding is used to make output size is same as the input size

p - padding amount = no. of rows/columns
that we will insert in top, bottom,
left and right of the image

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

Problems with convolution are

1) Shrinked output

2) Throwing away a lot of information
that are in the edges.

Soln: Pad the input image before
convolution by adding some rows
and columns.

strides :- It tell us the number of
pixels we will jump when we are
convolving filter

Input $n \times n$ * filter $f \times f$ p - padding
s - stride

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

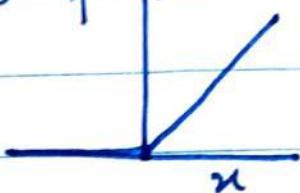
The convolution operation shrinks
matrix if $f > 1$

ReLU After every convolution operation, ReLU operation has been used.

Rectified Linear unit is a non linear activation function

$$\text{output} = \max(0, \text{Input}) \quad f(x)$$

$$f(x) = \max(0, x)$$



ReLU is an element wise operation, that measure if applied per pixel and replace all negative pixel values in the feature map by zero

- Convolution is a linear operation
non-linearity is introduced by using a non linear activation function

like ReLU

1	-2	8
-3	-1	1
1	1	-5

ReLU

1	0	8
0	0	1
1	1	0

feature -

CNN uses pooling layers to reduce the sized input to speed up computation and to make some of the feature it detect more robust.

steps - Pooling (spatial pooling) layer
Also called as subsampling or downsampling.

- It is applied after convolution and ReLU operations
- It reduces the dimensionality of each feature map but retains the most important information
- Since the number of hidden layers required to learn the complex relations present in the image would be large
- We apply pooling operation to reduce the input feature representation thereby reducing the computation power required for the network
- Spatial pooling can be of different types :
① Max pooling ② Average pooling ③ Sym pooling

① Max Pooling / subsampling

Once we obtain the feature map of the input, we will apply a filter of determined shape across the feature map to get the maximum value from that portion of the feature map.

- It is also known as subsampling because from the entire portions of the feature map covered by filtering kernel we are sampling one single maximum value

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

$\xrightarrow{\text{2x2 filter}}$
 $\text{stride} = 2$
 $p = 0$

4	7
8	5

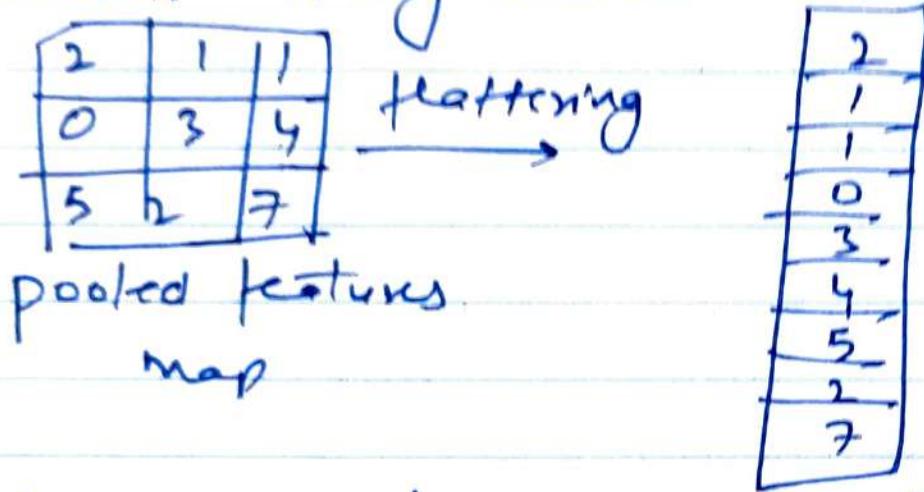
feature map 4x4

② Average pooling

- ③ Sum pooling compute the sum of all elements

Step 3 flattening

After a series of convolution and pooling operations on the feature representation of the image, we will flatten the output of the final pooling layer into a 2D single long (continuous) linear array vector.



The process of converting all the resultant 2D array into a vector called flattening.

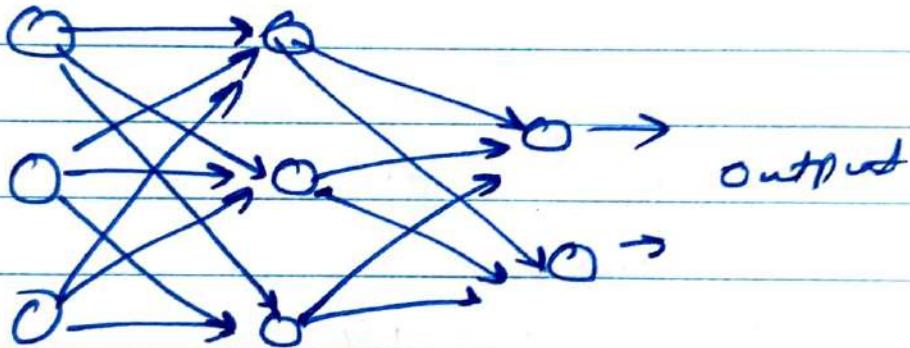
CNN + fc layer is too high level of information merge produced by convolution and pooling layer for classifying the input image into various classes based on the training dataset.

Fully connected means every neuron in the previous layer is connected

to every neuron in the next layer.
• sum of output probabilities from the FC layer is one.

Fully connected uses a softmax activation function in the output layers

The softmax function takes a vector of arbitrary real valued scores and transforms it to a vector



Step 7: Fully connected layer (dense layers)

Flatten output is feed as input to the fully connected ANN

- ANN have varying number of hidden layers to learn the non linear complexities present with the feature representation.

* sum of output probabilities from FC layer is one

Fully connected uses a softmax activation function in the output layer
→

the softmax function takes a vector of arbitrary transition real valued scores and transform it to a vector of values b/w zeros and one that sum to one

Working of CNN

Input Image

convolution (multiply filter)

+ ReLU
↓ feature map

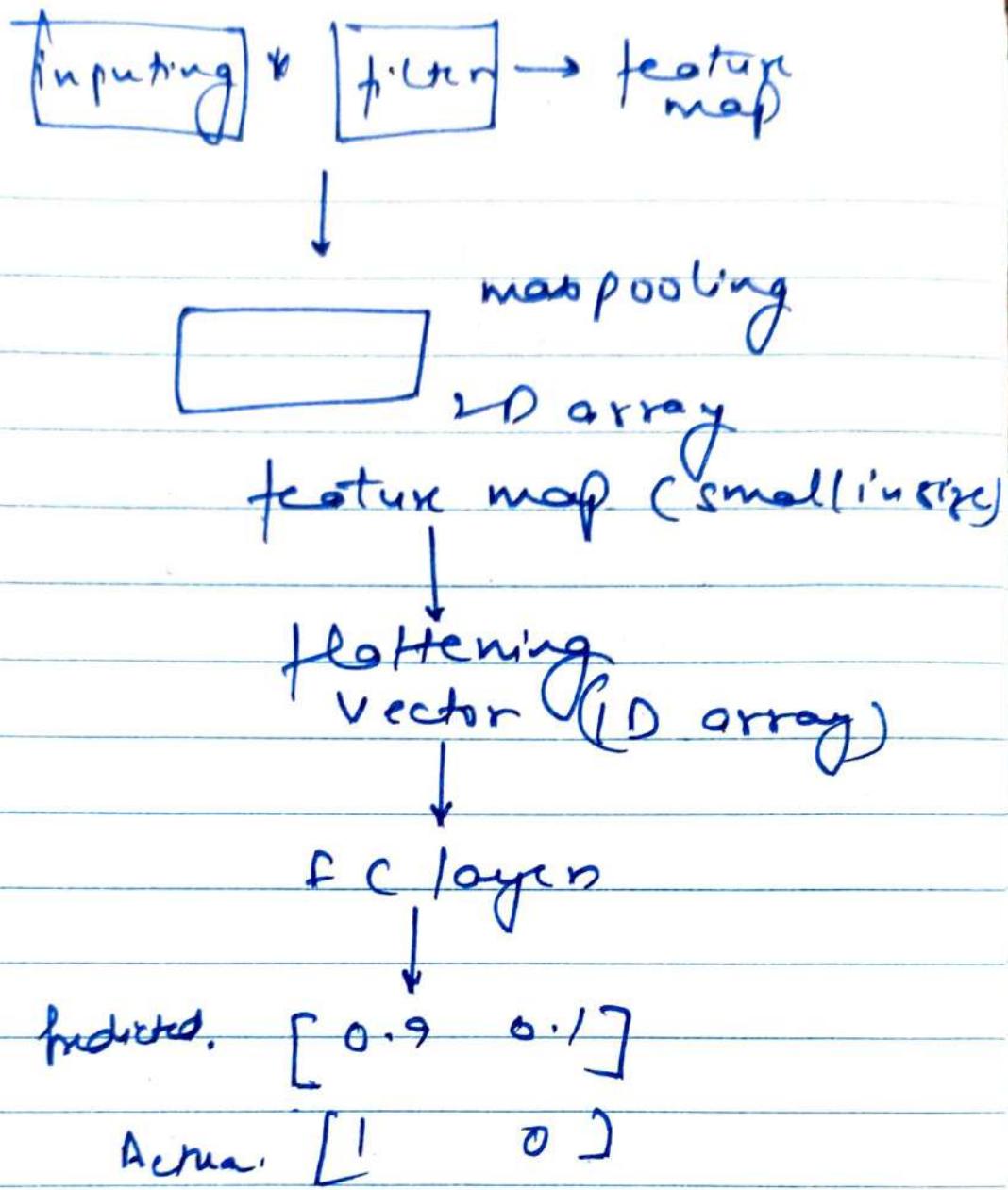
Max. pooling (reduction size of feature map)

Flattening

Dense / fully connected layer

↓
output vector

* One filter creates one corresponding feature map.



Training of CNN

- ① Initialize all filter weight with random values
- ② Forward propagation C / Input \rightarrow convolution \rightarrow ReLU \rightarrow pooling \rightarrow fc \rightarrow o/p
 probability to recognize
- ③ Calculated the total errors (SSE)

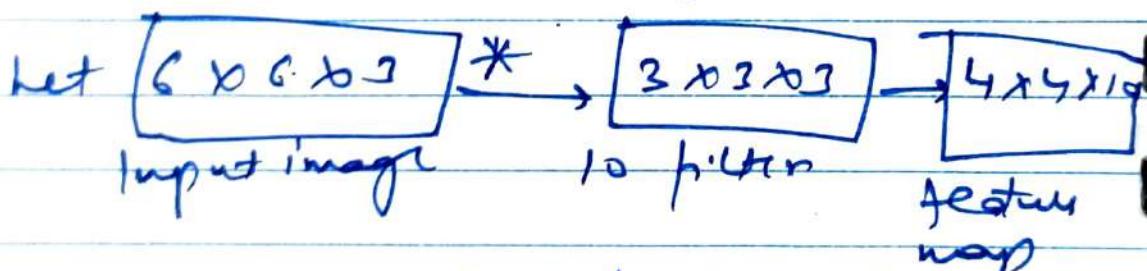
$$\frac{1}{2} \sum (target_p - output_p)^2$$

4. Backpropagation (using gradient descent) update weights and parameters. Repeat till minimum loss/error. Repeat for all input images (training set)

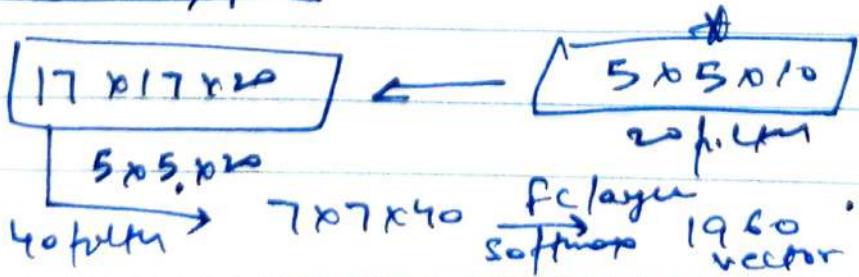
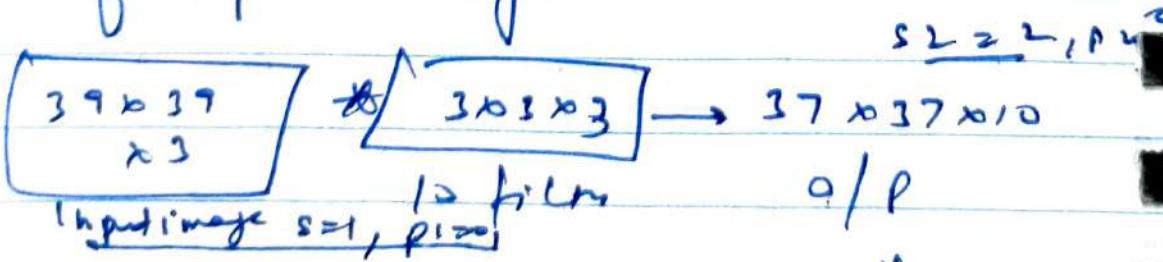
Convolution over 3D images

We will convolve an image of height, width, number of channel with a filter of a height, width, same number of channels.

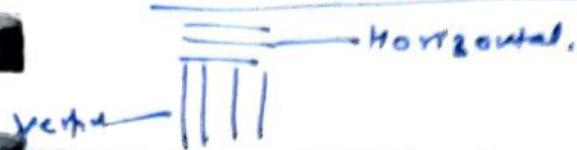
No. of channels in image = no. of channels in filter



We use multiple feature filters to detect multiple features or edges of input image

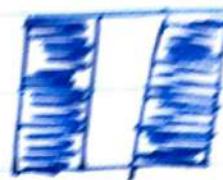
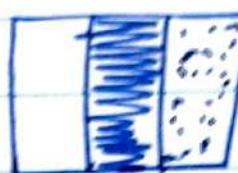


Vertical and Horizontal edge detection



Vertical Edge detection

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



Horizontal Edge detection

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

1x1 Convolutions (Network in Network)

Problems of CNN

3x3 filter will be applied
in 3x3x3 block

3x3x3 blocks to create
a single value

6x6x3

3x3x3

4x4x2 filter

- If the size of filter is relatively large like 5x5 or 7x7 pixels, it can result considerably more parameters (weights) which require large computation to

perform the convolutional operations.

- Convolutional and pooling layer reduces the height and width of feature map leaving the depth number of channels or number of filters unchanged (number of feature maps)

- So deep CNNs requires a layer that can down sample or reduce the depth or number of feature maps.

3×3 filter will be applied on $3 \times 3 \times 64$ blocks to create a single value to make up the single output feature map.

When the depth of network increases, the depth of input or number of filter used in convolutional layers often decreases, resulting in an increase in number of resulting feature maps

- A large map number of feature maps is a concern can cause a problem as a convolutional operation must be performed down through the depth of the input

Solution Use a 1×1 filter to downsample or number of feature maps.

1×1 convolution can be used to reduce only the depth as it slices through the input volume with just 1×1 unit.

- It save computation task
- It adds non linearity
- It allows to learn more complex function.
- It generates an element wise product the depth, slices through the entire volume.
- Useful inner larger filter sizes and used such as 5×5 , 7×7
- It is used for dimensionality reduction
- It is so useful in many CNN models

$6 \times 6 \times 3^2$ * $1 \times 1 \times 3^2$ Input #filter 5

$6 \times 6 \times 5$ Output

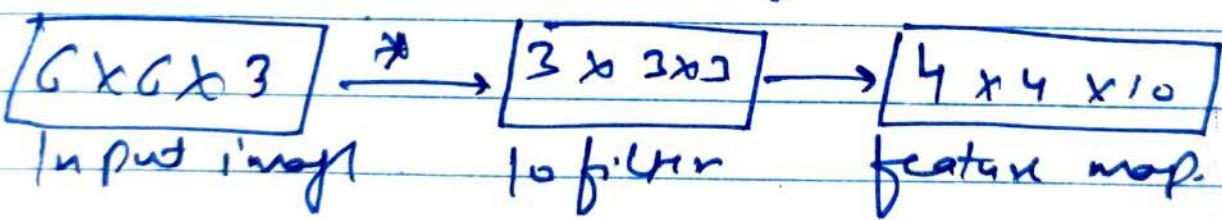
↓
- by channel

- It has been used in a lot of modern CNN implementations like ResNet & inception models
- It is useful when we want to shrink the number of channel called feature transformation. This save a lot of computations.

Convolution over 3D images

We will convolve an image of height, width, number of channels with a filter of a height width, same number of channels.

no. of channels in image = no. of channels in filter



We use multiple features to detect multiple features or edges of input images

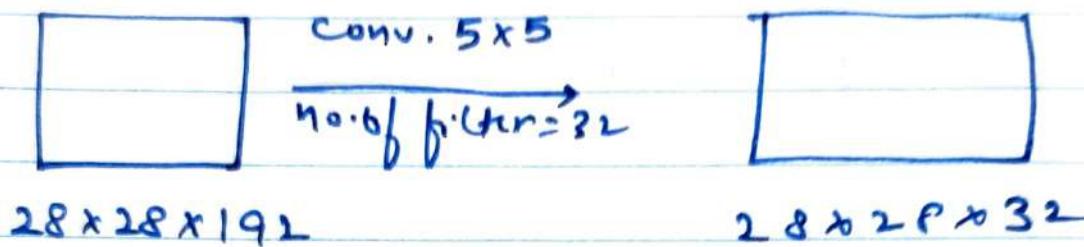
Inception network During design of a CNN you have to decide all the layers yourself

Choices —

3x3
5x5

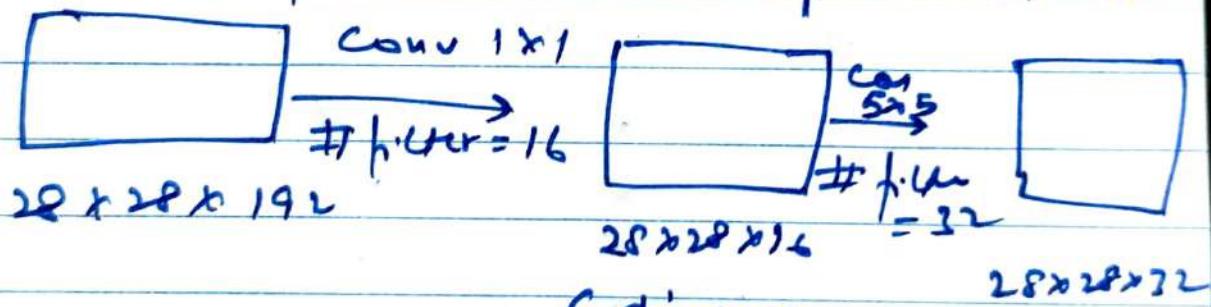
max pooling layer

5x5 convolutions



$$\begin{aligned} \text{computation cost } & 28 \times 28 \times 5 \times 5 \times 192 \times 32 \\ & = 120 \text{ M} \quad (\text{high computation cost}) \end{aligned}$$

Using a 1×1 convolution with 5×5 convolution to reduce computation cost



(dimensions reduction)

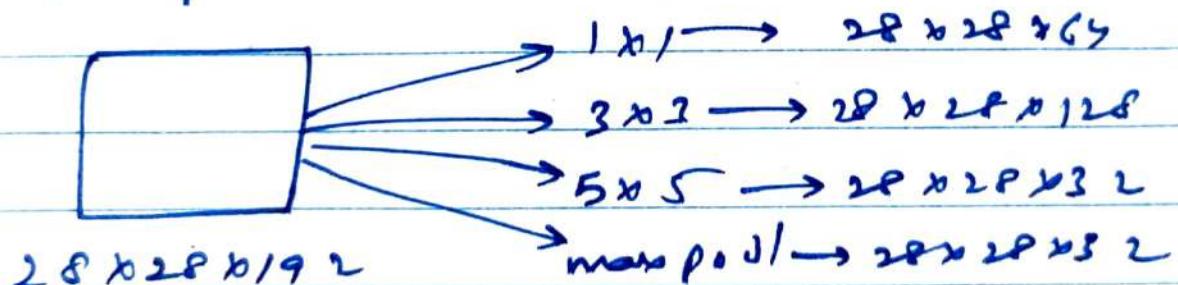
$$\begin{aligned} \text{— computation cost } & (28 \times 28 \times 16) \times (1 \times 1 \times 192) \\ & + (28 \times 28 \times 16) \times (5 \times 5 \times 16) \\ & = 12.4 \text{ M} \quad (\text{approx}) \end{aligned}$$

— Hence 12.4 M is immensely smaller than 120 M ~~IF~~ it means 1×1 convolution can help to reduce model size which is they also help to reduce computation

Inception network, instead of choosing a $1 \times 3 \times 3$ or 5×5 filter or a pooling layer we apply all together

- This makes the network architecture more complicated but remarkably improved performance as well
- The inception network or inception layer says - instead of choosing the desired filter size in a conv. layer or whether we want a convolutional layer or a pooling layer lets apply all of them

Inception module, naive version.



$$= 64 + 128 + 32 + 32 = 256$$

$$28 \times 28 \times 192 \rightarrow 28 \times 28 \times 256$$

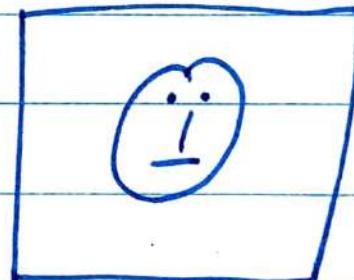
Input

Output

High computation cost

Here we used all the neurons & Pool we might want and will let the neural network learn and decide which it wants to use most

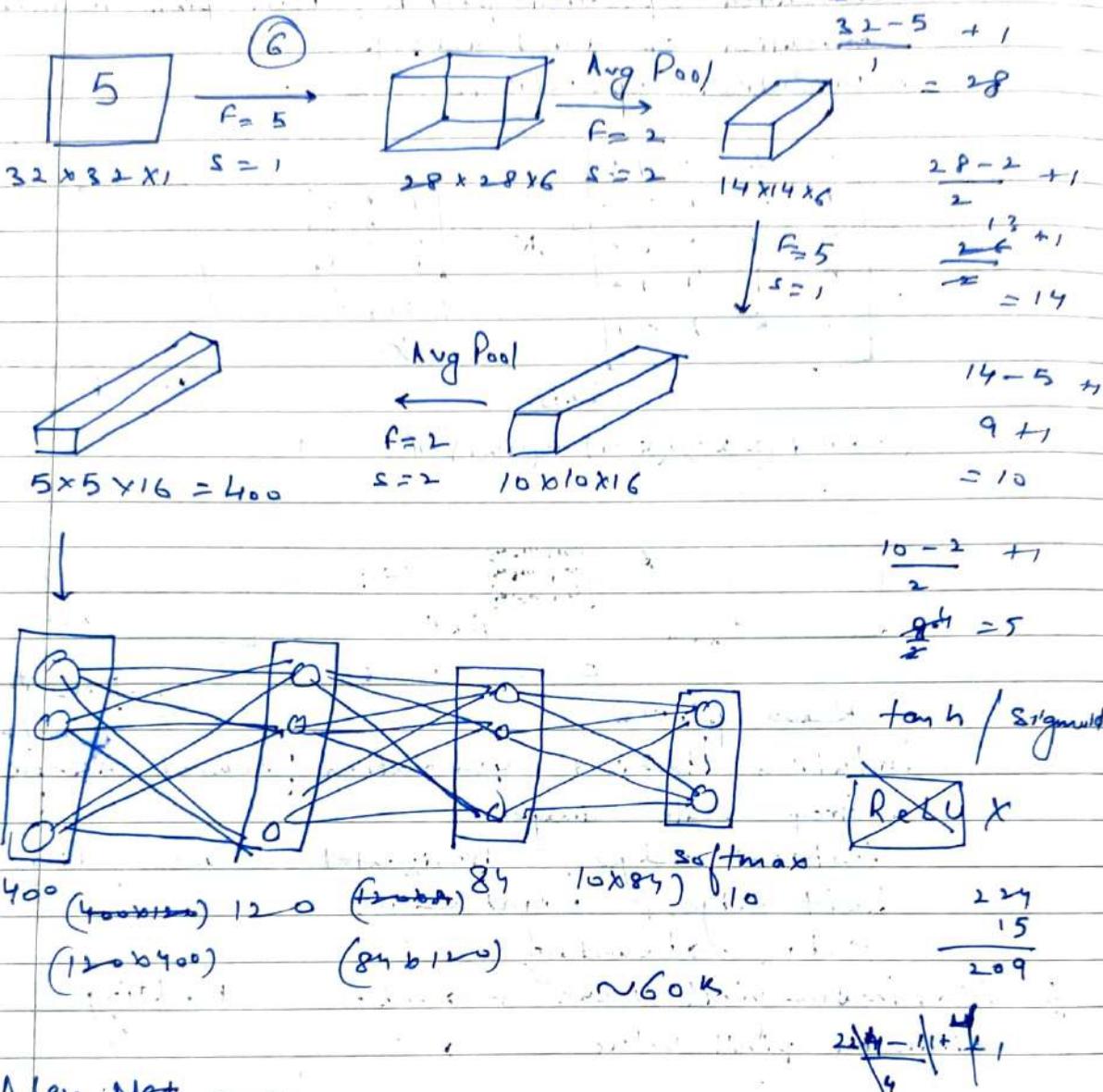
Inception layer allows the internal layers to pick & choose which filters' size will be relevant to learn the required information so even if the size of the face is in the image is different, the layer works according to recognize the face



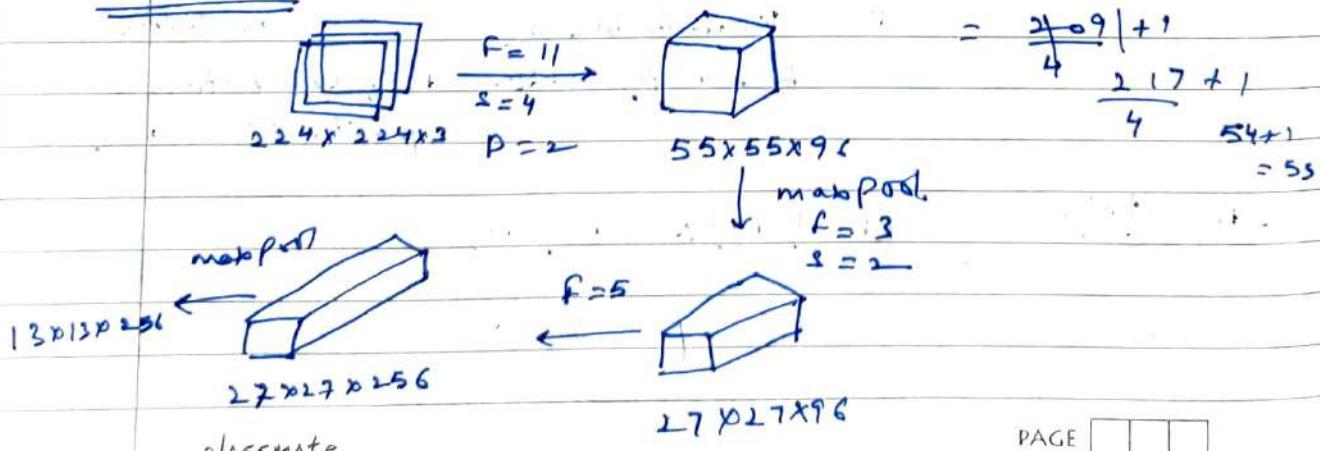
Sparse connections
Parameters sharing

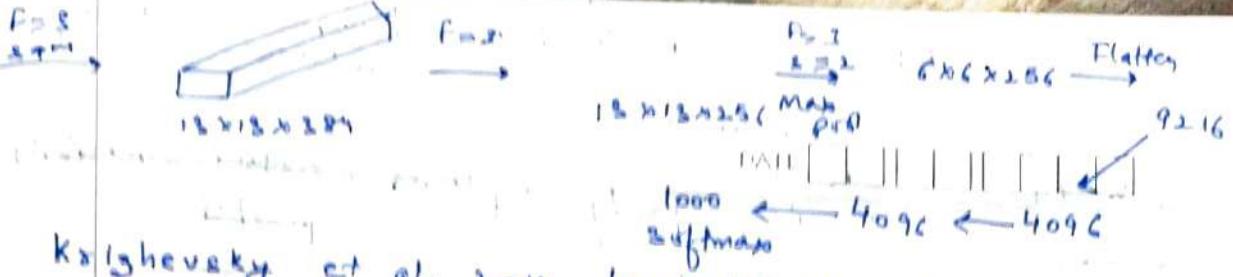
LeNet - 5 architecture:

LeCun et al., 1998, Gradient-based learning
Applied to Document Recognition



Alex Net

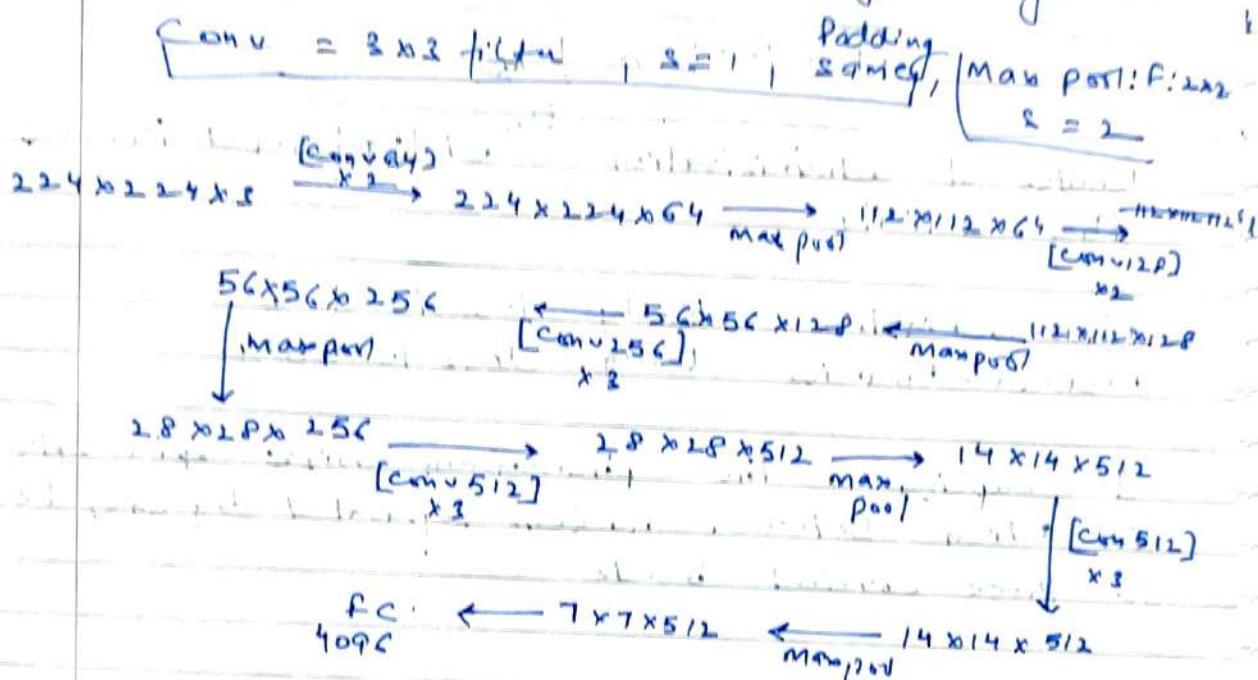




VGG-16

Simonyan & Zisserman (2015), Very Deep Convolutional Networks for Large-scale Image Recognition

ReLU used



16. Represent no. of learning parameters
~138 Million parameters

VGG-19

ResNet Architecture. Kaiming He, Xiangyu Zhang, Ren, & Jian Sun "Deep Residual Learning for image Recognition" CVPR 2016.

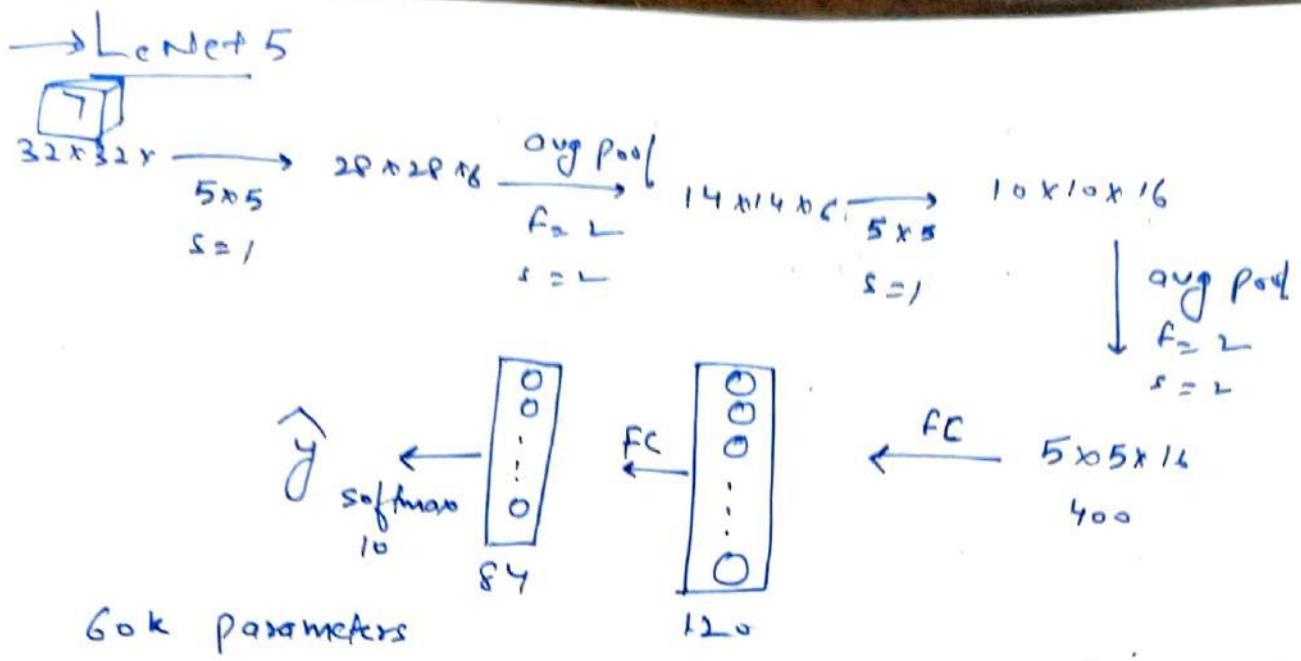
Dotted line: the dimension

Skip connection

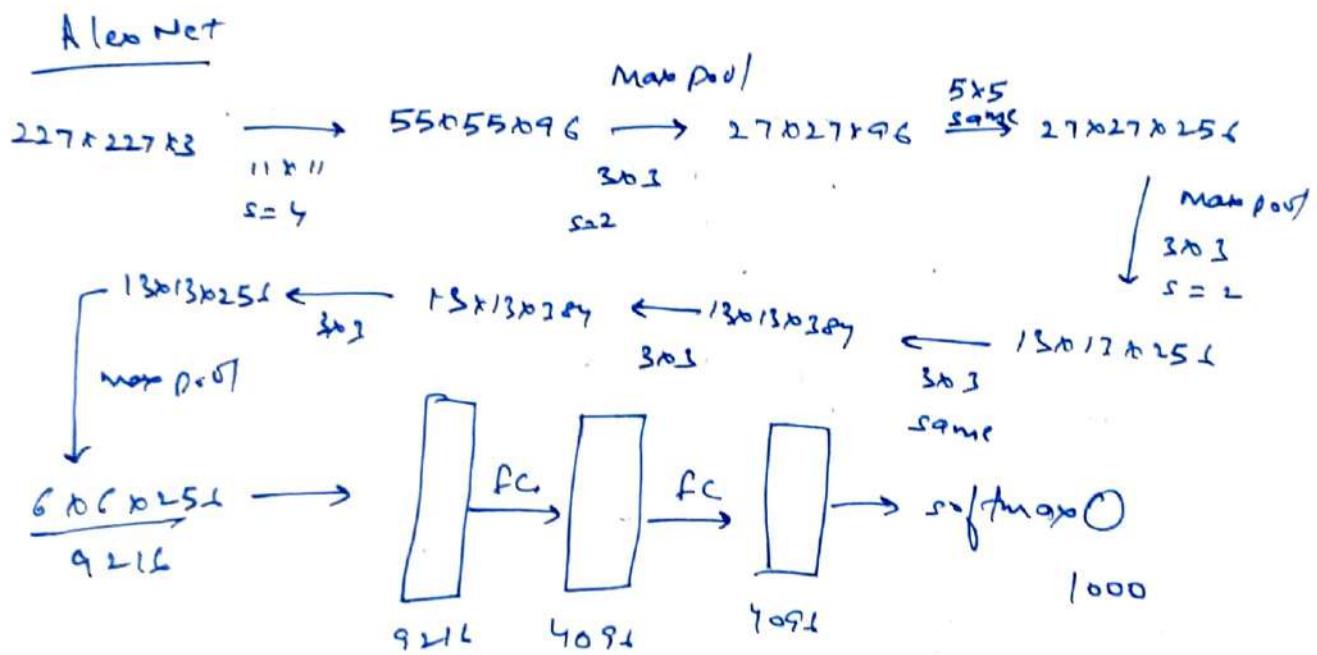
152 Layers, Less overfit

Residual block

classmate



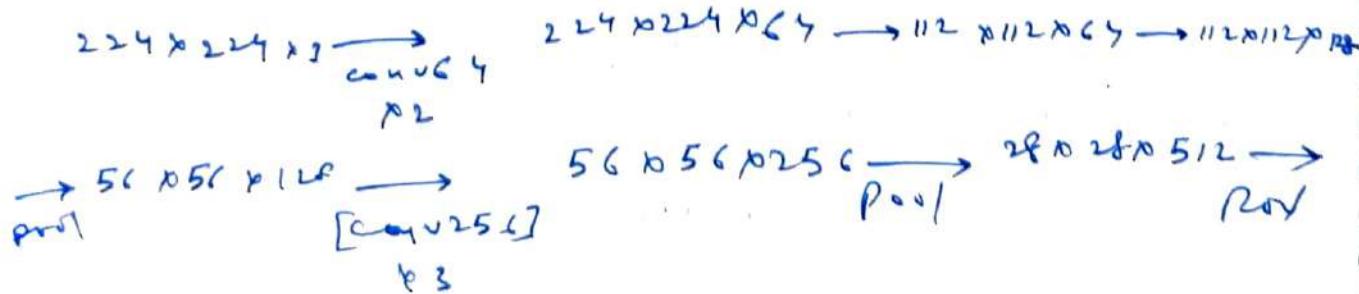
Advance comments! . Sigmoid / tanh .

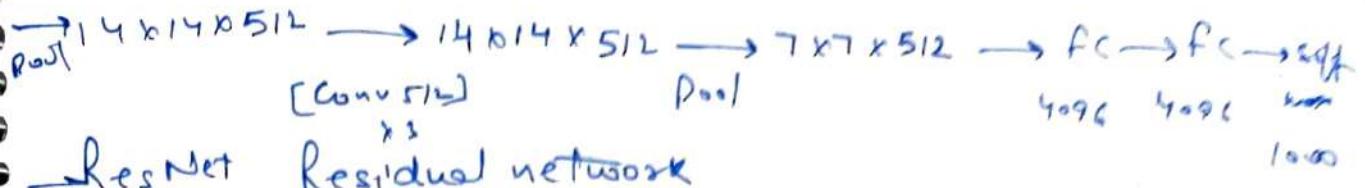


VGG-16

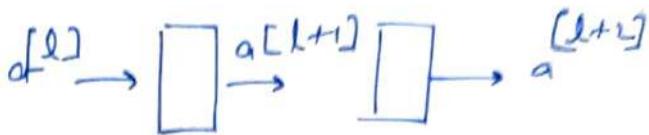
Conv = 3x3 filter, $s=1$, same

Max pool = 2x2, $s=2$





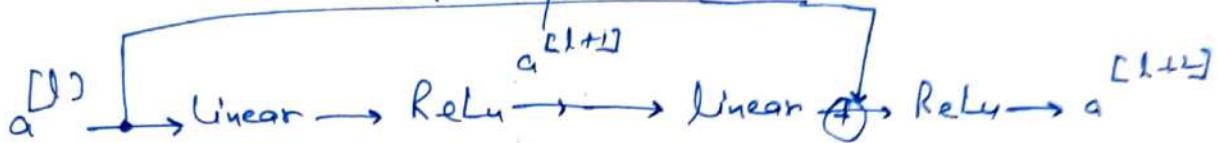
Residual block



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

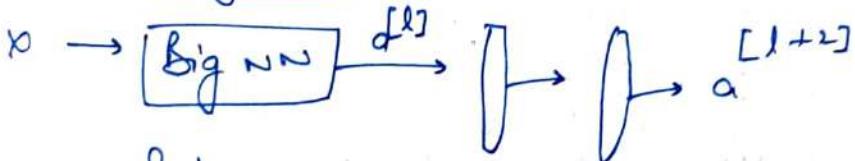
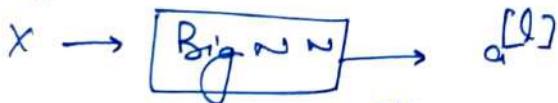
$$z^{[l+2]} = w^{[l+2]} a^{[l+1]} + b^{[l+2]} \quad a^{[l+2]} = g(z^{[l+2]})$$

shortcut / skip connection



why ResNets work

why do residual networks work?



ReLU

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

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

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

— Network in network (1x1 convolution do)

— Inception network

→ Using open source implementation

→ Data augmentation.

- Mirroring
- Random cropping
- Rotation
- Shearings
- Local warping ...

color shifting

Advanced:

PCT

PCT color Augmentation.

Implementing distortions during training

→ state of computer vision

Data V.s hand engineering

→ what Object localization

what are localization and detection?

① Image classification, classification with
localization, Detection

b_x, b_y, b_h, b_w

- Defining the target / label y_j

e.g

1	if object is present
0	if no object
1	if object
0	if no object
1	if object
0	if no object

$$y = \begin{cases} P_C \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{cases}$$

is there any object

$$\mathcal{L}(g, y) = (g_1 - y_1)^2 + (g_2 - y_2)^2 + \dots + (g_f - y_f)^2$$

- Landmark detection
- Sliding windows detection
- Convolutional Implementation [Sliding windows
Sermanet et al., 2011 Overfeat! Integration, localization,
and detection using convolution networks]
- Bounding Box Predictions
 - Yolo algorithm.
 - Labels for training
 - for each grid cell
 - label for training
 - For each gridcell:
 - assign mid point of gridcells
 - relative to grid cell
- Intersection over Union
 - Intersection over union
 - = $\frac{\text{size of intersection}}{\text{size of union}}$
 - correct iou if $\text{iou} \geq 0.5$
 - You measure of the overlap b/w two bounding boxes
 - Non max suppression
 - Discard all boxes with $P_c \leq 0.6$
 - Anchor boxes
 - for multiple objects
 - predefined two shapes: anchor boxes (L)
 - Yolo algorithm
 - Region proposal R-CNN
 - Propose regions. classify proposed regions
one at a time. Output label + bounding box
 - fast R-CNN
 - faster R-CNN