

Biological Inspiration : Visual Cortex (4.1) :-

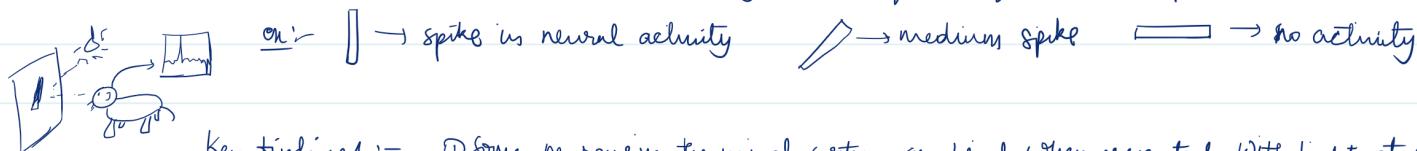
aka CNN / ConvNets.

→ Mainly used for visual tasks such as object recognitions etc,

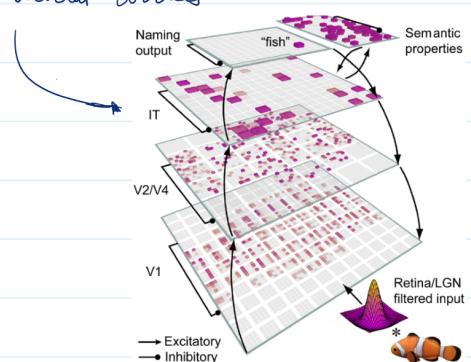
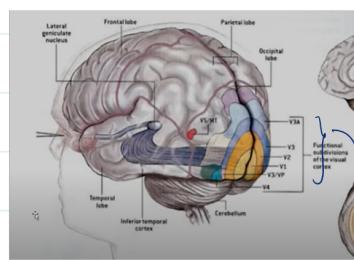
→ Breakthroughs in (DL + CV) happened in the last 10 years.
 ↗ Computer Vision

→ Hubel & Wiesel experiments in 1950-80s ~

on animals. certain neurons had more activity when light edge is of a shape -



- Key findings :-
- ① Some neurons in the visual cortex are fired when presented with light at a specific angle. ex:- V1 : primary visual cortex → detects edges.
 - ② There are neurons for edge detections, motion, depth, color, faces, shapes etc.,
 - ③ There is a hierarchical structure among visual cortices.



→ Lots of image processing & CV is inspired by Hubel & Wiesel research.

→ ConvNets take advantage of this research & biological structures.

Convolution : Edge detection on image (4-2) :-

→ performed by Primary Visual cortex (V1)

ex:-

$$\begin{matrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 255 & 255 & 255 & 255 & 255 \\ 255 & 255 & 255 & 255 & 255 \end{matrix}$$

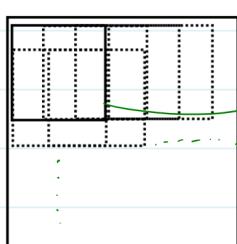
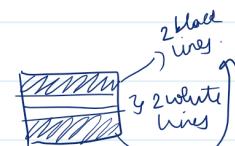
* not multiplication convolution.

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

$$\begin{matrix} 0 & 0 & 0 & 0 \\ -120 & -100 & -100 & -100 \\ -100 & -100 & 400 & 100 \\ 0 & 0 & 0 & 0 \end{matrix}$$

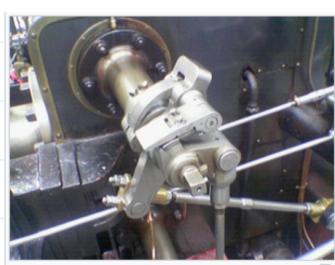
max max norm

$$\begin{matrix} 255 & 255 & 255 & 255 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 155 & 255 & 255 & 255 \end{matrix}$$



→ Sobel vertical edge detector :

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



$$G_x = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix} \quad G_y = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

Adjoint matrices obtained from both for both horizontal & vertical

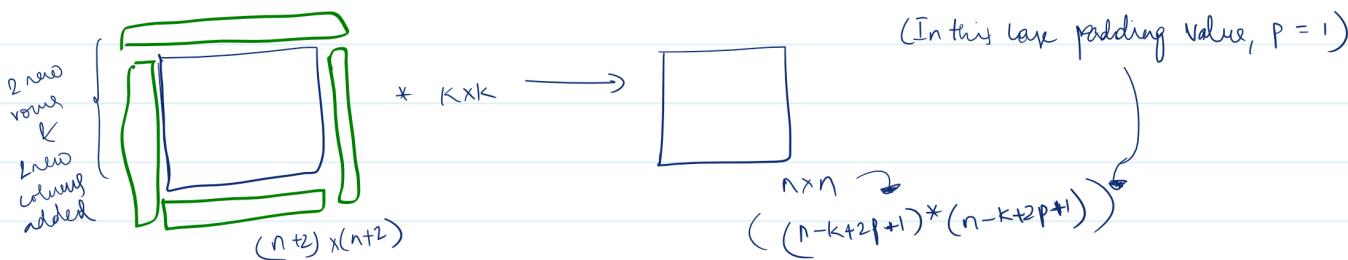
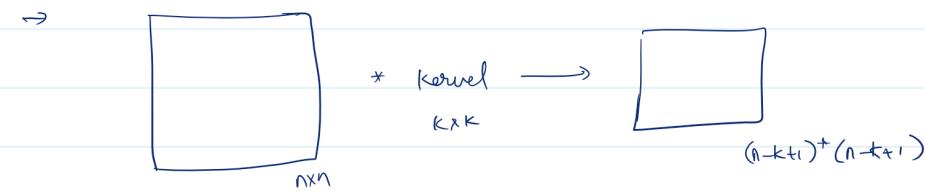
→ Kernel size is not fixed. Not always 3×3

→ Kernel value can also change.

→ Box filter = $\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ → blurs the image when sum is divided by number of coeff.

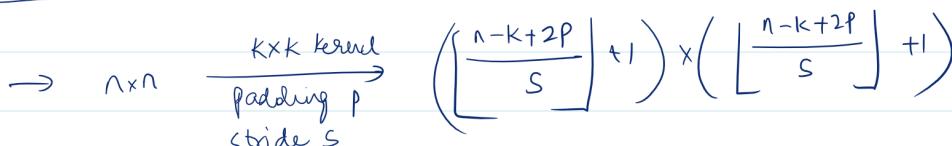
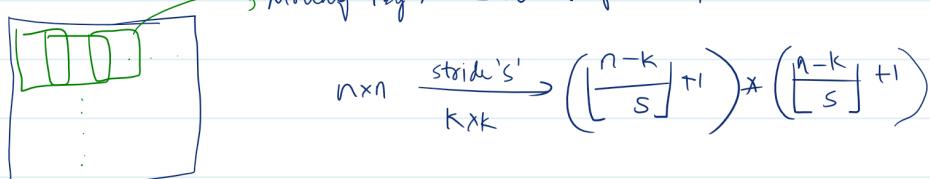
$$\begin{bmatrix} \frac{1}{16} & \frac{1}{8} & \frac{1}{16} \\ \frac{1}{8} & \frac{1}{4} & \frac{1}{8} \\ \frac{1}{16} & \frac{1}{8} & \frac{1}{16} \end{bmatrix} \rightarrow \text{blurs the image directly.}$$

Convolutions : Padding & Strides (4.3) :-



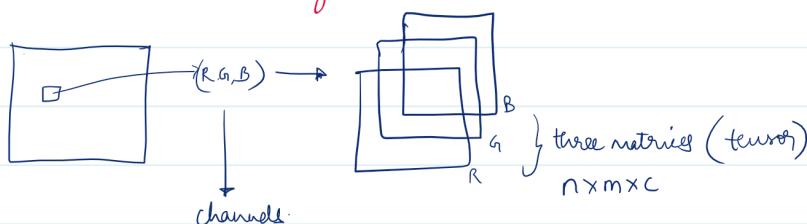
→ Padding can be filled by (i) zeros : Most common. Adds more edges to the image though. These edges are not learned again in backpropagation.
(ii) Same Value Padding : Not common. Causes conflicts in case of multiple neighbors.

Stride :-

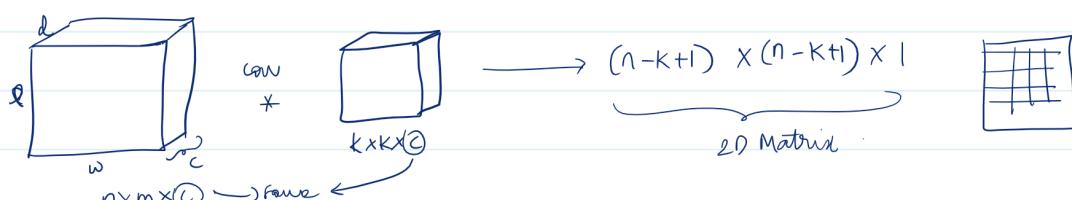


→ Padding at bigger output, stride for smaller output.

Convolution over RGB Images (4.4) :-

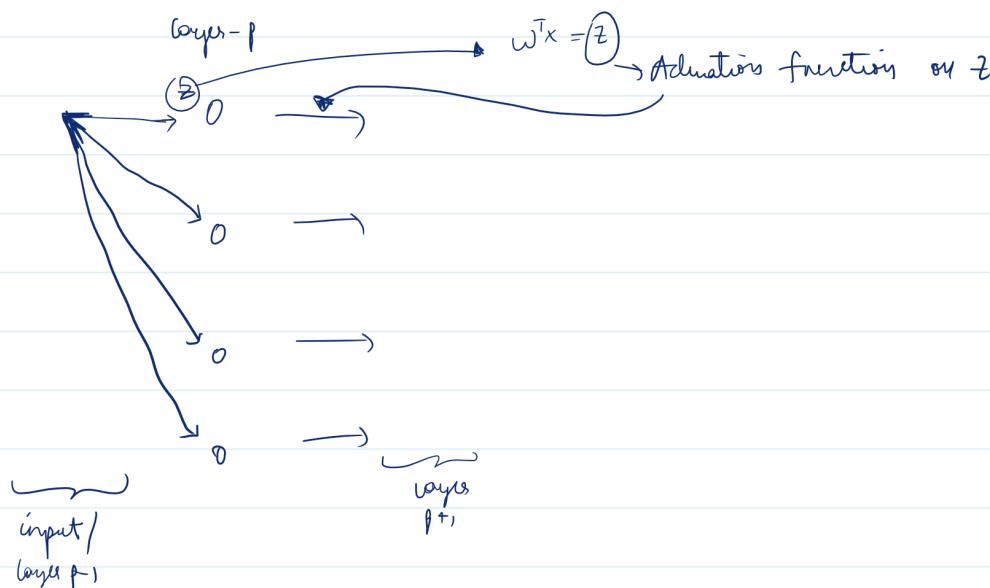


→ Convolution performed same as in 2D i.e., component wise multiplication and addition



Convolutions layer = (4.5)

→ In normal MLP,



→ But in CNNs,

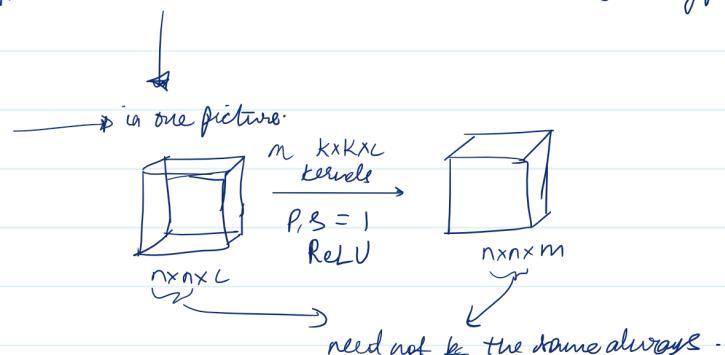
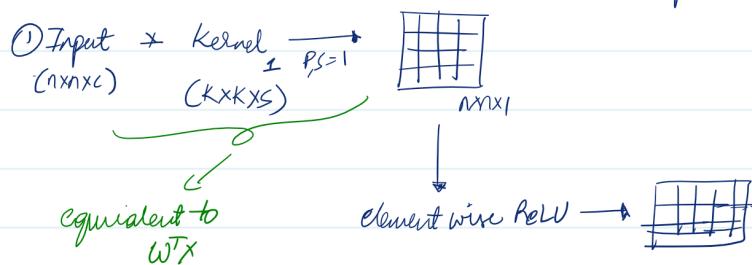
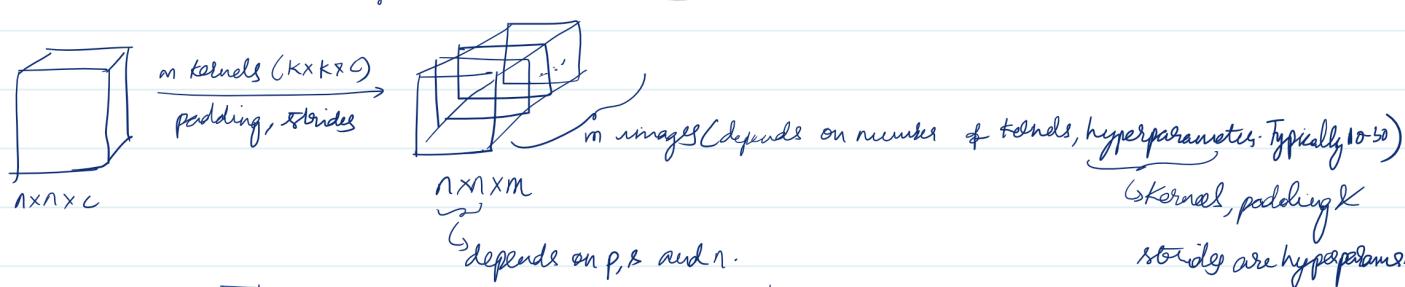
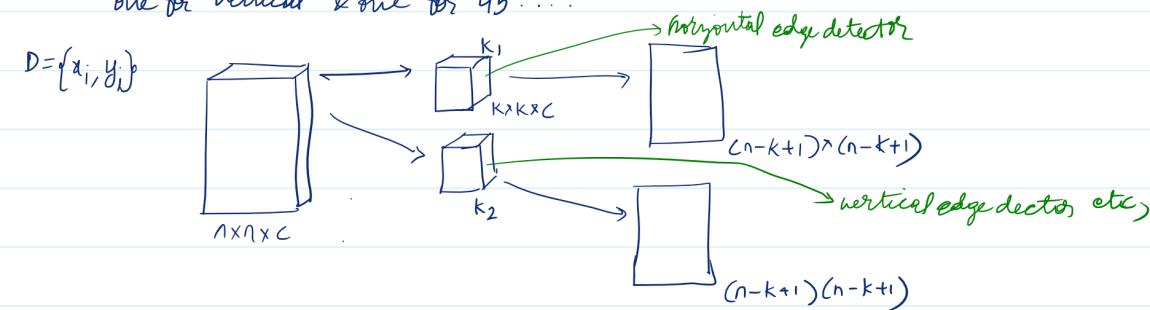
multiple edge detectors \Rightarrow multiple kernels
(1 edge detector = 1 kernel)

So in CNN, we learn the kernel matrices in backpropagation

MLP: learn weights

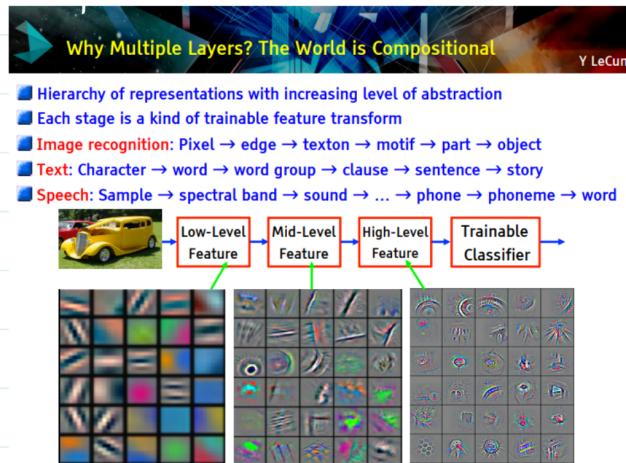
CNN: learn kernel matrices. ($I_{K \times K \times C}$)

for ex:- instead of hardcoding one kernel to identify edges, we could have one kernel for horizontal edge, one for vertical & one for 45° ...



→ In a deep multi-layer CNN, works like human visual cortex.

first layer identifies parity, second edge, third texture, ..., part, object etc.)



→ Using ReLU is not mandatory in CNNs.

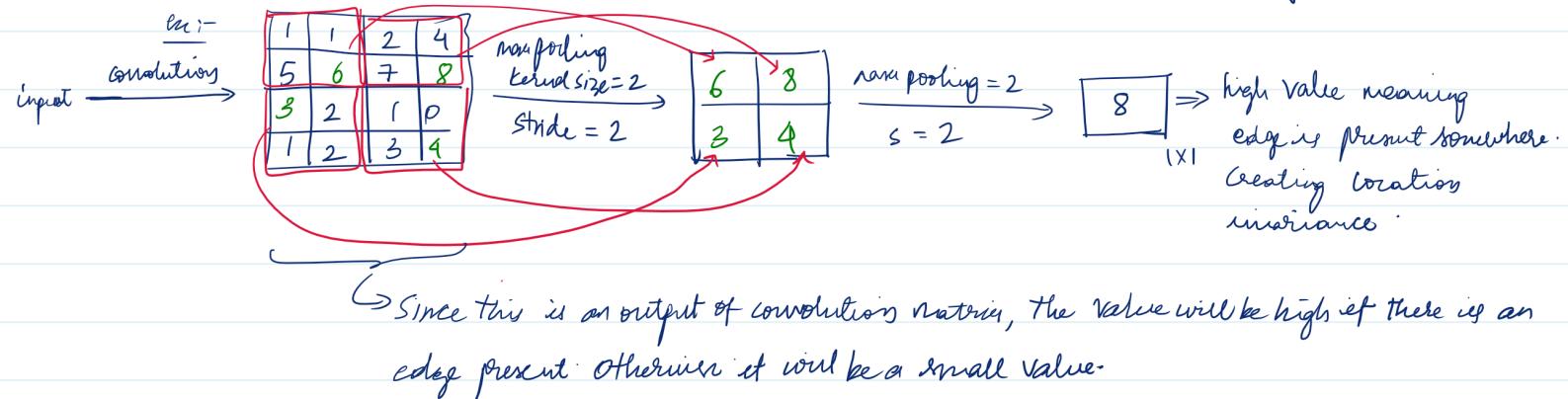
→ Kernels are initialized in the same way as in MLP.

Max Pooling (4.6) :-

→ Location invariance :- We want model to detect object irrespective of its location in image. Similarly scale invariance, rotation invariance etc.

In human body at higher levels of visual cortex, (V7 etc) neurons are fired as soon as object is recognised irrespective of scale, size etc. Max pooling tries to replicate this.

→ Max pooling can also be thought of as a kernel. Finds the max value in the kernel shape & size.



→ Very popular in ConvNets. Avg-Mean pooling is also gaining popularity.

→ Capsule Network captures both presence of an object & its local data.

If eyes, ears, mouth & nose positions are jumbled, CNN will recognise a face, but CapsNet won't.

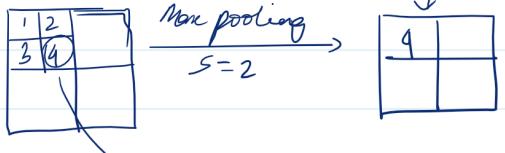
→ Repeated kernels automatically get rid of scale invariance. Data augmentation takes care of rotation invariance.

CNN Training : Optimisation (4.7) :-

→ In order to calculate derivatives, operations need to be differentiable.

CNN operations :- Convolution layer :- Convolution + ReLU → differentiable
Max Pooling

→ dot product + addition intuitively. In MLP, its element wise operations. In Convnets, its pixel wise operations.



derivative of g is only passed to x in conv matrix.

$$\text{ie, } \text{max}(x) = \frac{\partial g}{\partial x_i} = \begin{cases} 1 & \text{if } x_i = \text{max}(x) \\ 0 & \text{otherwise} \end{cases}$$

derivative is passed only to this value. This is done because no matter what happens to other vals, it doesn't affect max pool. Thusly the end result -

→ Same thing happens to mean pool/avg pool.

→ We can use dropouts in CNN as well just like in MLPs.

Read first comment links for more content.

→ if we invert kernel matrix in convolution by 180°, it becomes convolution operation

ie,

a	b
c	d

 \Rightarrow

d	c
b	a

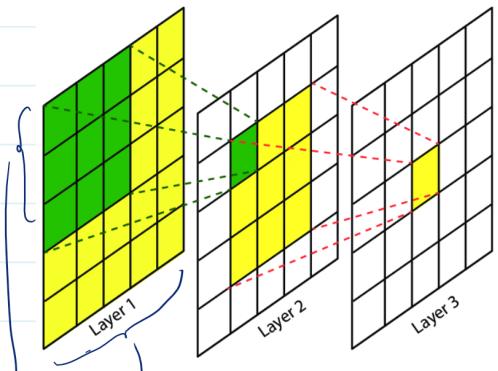
↙ ↘
convolution.

→ Max pooling selects the brightest pixels from the image. Useful when image is dark & we are interested in only brightest pixels of image. Mean pooling smoothes the image & sharp features may not be extracted.

Receptive Field and Effective Receptive Field (RF) :-

Receptive Field :- The regions on image on which kernel operates at a given time 'T'.

Effective Receptive Field :- In deep CNN, the region on original image on which kernel at a deeper layer indirectly works on.



• Receptive field for Layer 2.

↳ effective receptive field for Layer 3.

ImageNet (4.10) :-

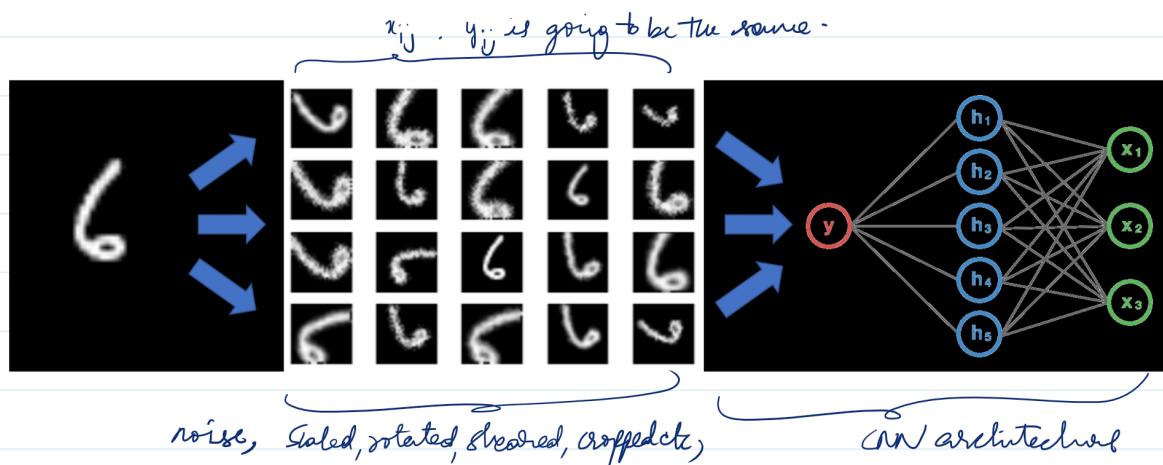
→ CNN used AlexNet in 2012 gave very good result.



→ AlexNet vs LeNet.

Data Augmentation (4.11):-

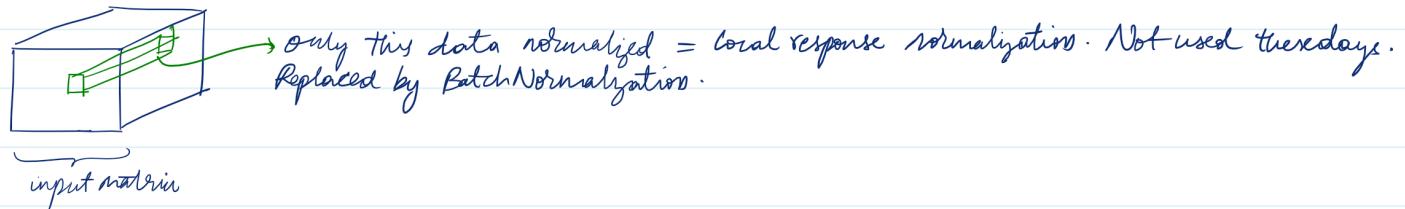
→ We want the models to be robust & work even when image is rotated, scaled etc.



→ Why? We can get invariance like Scale invariance, rotation invariance etc. Also can increase the size of the dataset.

→ Augmentation prevents model from overfitting by introducing diversity to the dataset.

AlexNet (4.13):-



→ AlexNet used dropouts, ReLUs & multiple GPUs to train. Gave good score.

→ Objective of dense layer is to take results from convolution/pooling and use them to classify image into a label.

VGG Net (4.14):-

→ By Andrew Zisserman (key figure in CV). AlexNet uses all sorts of kernel sizes. VGG tries to use simple fixed kernel sizes i.e., 3x3 kernels, $S=1$, $p='same'$ and Maxpool $S=2$, 2×2 .

→ VGG = Visual Geometry Group.

→ 2 types VGG-16 & VGG-19. Default networks used these days.

16 layered & 19 layered network.

→ Pretrained available in Keras.

→ In training, weight decay is regularization = 5×10^{-4} & dropout ratio = 0.5.

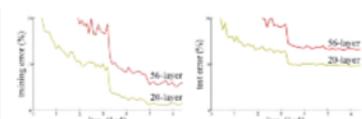
→ Regularization can be applied at all layers including conv layers because they contain parameters too.

→ In AlexNet MaxPooling done after almost every conv layers. In VGG, it's done after 2-3 conv layers.

ResNet (4.15):-

→ Residual Networks. by Microsoft.

← Even after adding lots of layers, error not reducing even after dropout, lots of its etc.)



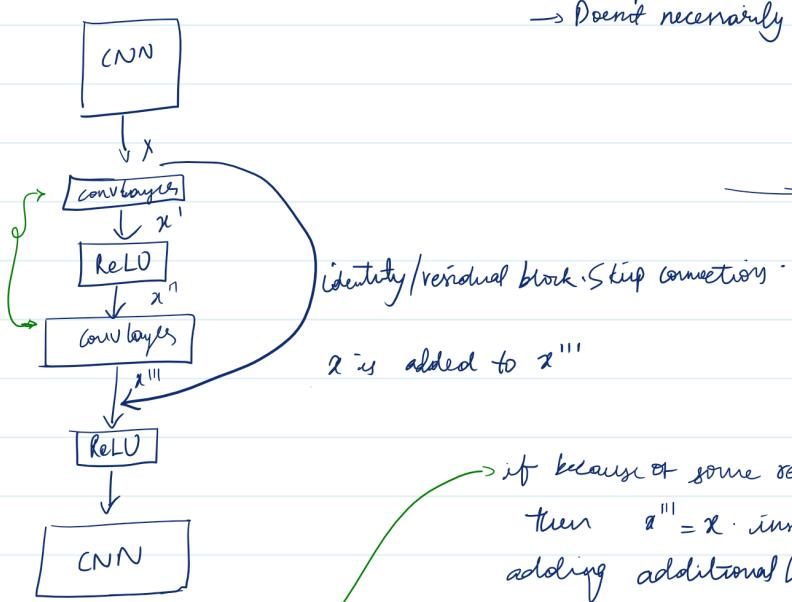
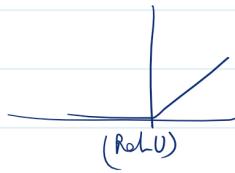
→ Doesn't necessarily have to be conv layers. Can be dot prod (MLP) too.

If x is +ve, $\text{ReLU}(x) = x$

$$\text{ReLU}(\text{ReLU}(x)) = x = \text{ReLU}(x)$$

If x is -ve, $\text{ReLU}(x) = 0$

$$\text{ReLU}(\text{ReLU}(x)) = 0 = \text{ReLU}(0)$$

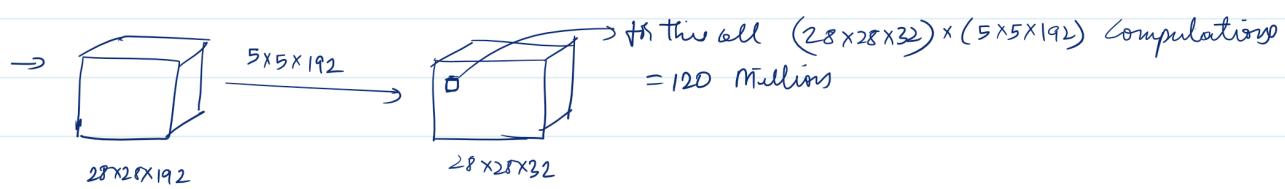
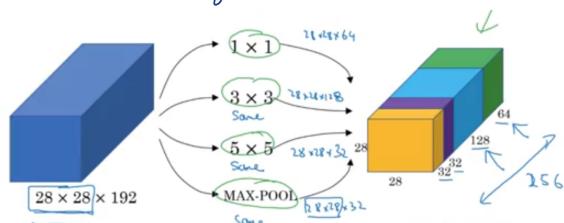


→ if because of some regularization, there 2 become 0 (unless), then $x''' = x$. instead of 0. ie, output will not be impacted & adding additional layers would not hurt performance - will skip over them.

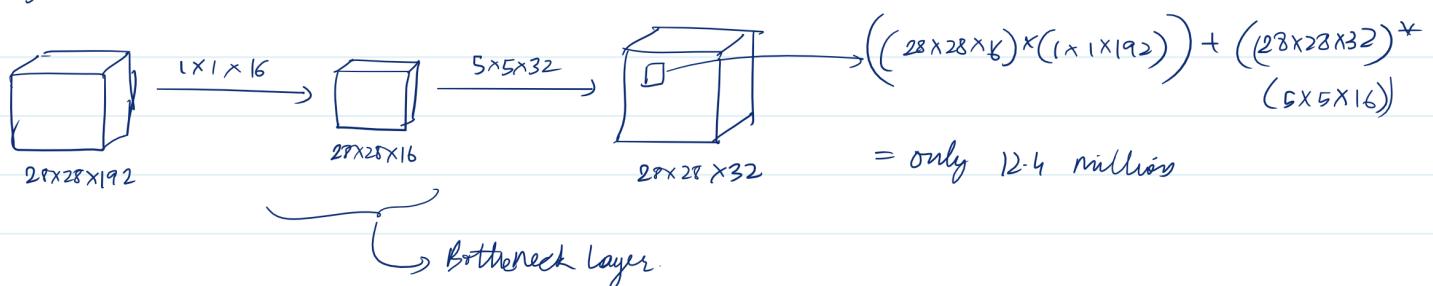
→ Must be before ReLU & x & x''' must be same size i.e., padding should be taken care of in conv layers.

Inception Network (4-16):

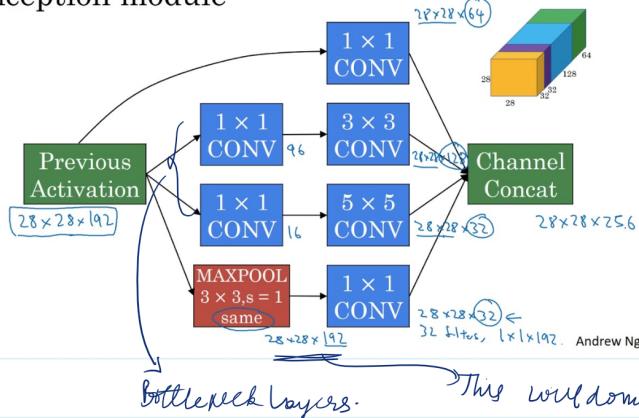
→ Instead of using just one kernel, use multiple kernels at once.



But instead



Inception module



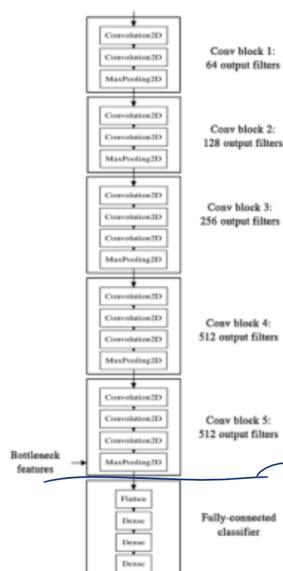
This will dominate the conv kernels above it. So we use 1x1 to reduce 192 to 32.

Transfer Learning (4.17) :-

→ Instead of building NN classifiers from scratch, we can reuse existing models, trained on diff dataset (like VGG16). We "transfer" learnings from one model & use it for our use case.

Case ① :-

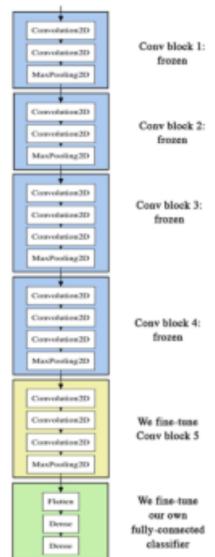
Here's what the VGG16 architecture looks like:



we only use these layers

Bottleneck layer - By this layer we have a vector representation of each image
ie, we use VGG as a sort of featureengineering tool till flatten if classification problem or till bottleneck it Segmentation problem (goal is to divide image into segments).

Case ② :-



we "freeze" these layers as they are anyway needed

We modify these layers acc to our needs using our dataset D. We should ensure learning rate is small so things don't change drastically.

Case ③ :- VGG16 + Imagenet to initialize new model & tune model entirely using our dataset D. Learning rate small

Case ④ :- retrain model from scratch. (Bad & expensive)

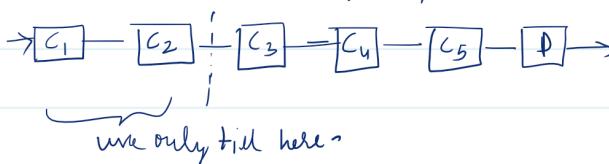
→ which case to choose depends on size of our dataset

Case ⑤ :- When $D \ll$ Imagenet & $|D|$ is small

Case ⑥ :- $|D| \approx$ Imagenet & $|D|$ is large: finetune the NN.

Case ⑦ :- $|D| =$ medium $D \ll$ Imagenet: finetune last layers.

Case ⑧ :- $|D| =$ small & $D \not\approx$ Imagenet



we only till here

Case ⑨ :- $|D| =$ large & $D \not\approx$ Imagenet, init model using VGG16 & retrain.

Fine tune = starting from a trained network, retraining on new dataset using very small update changes.

General trend in designing NN architecture

CNN → pool → CNN → pool → Dense → Dense (sigmoid/softmax)

Can replace with any number of dense layers as long as it makes sense.

Can also be refined as CNN - BN - Activation - Dropout - pool → . . . → Dense → Output.

Face Recognition from scratch in CNN:-

- 1000 images per class is good for CNN (if doing without transfer learning)
- last couple of layers with 50-100 images per person. Data Augmentation as well.
(Transfer learning) ↗ of vgg. ↗ Realistic setting (glasses, lighting etc.) and various angles.
- Good loss would be cross entropy & softmax last layer.
- If no time, use cloud APIs that are pre trained. (ex:- MS Advanced Facial Recognition Systems)