

DLP - Week - 9

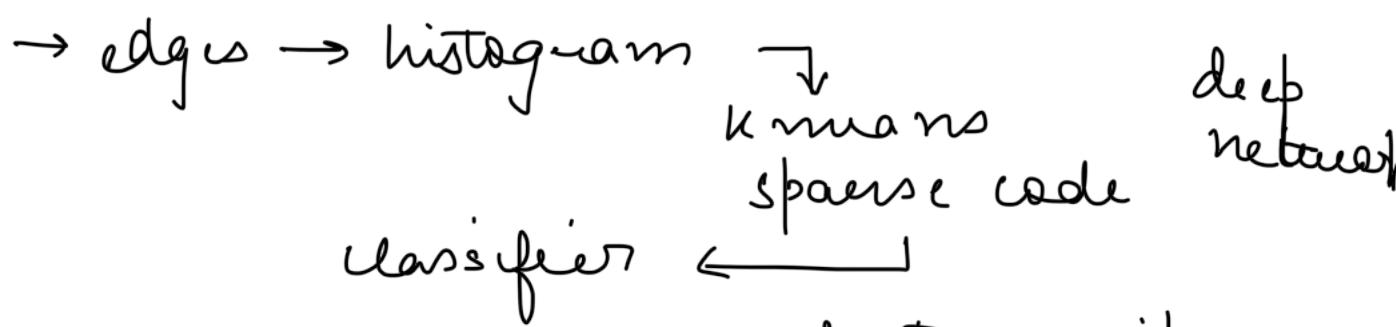
L1 Introducⁿ to Computer Vision

- we want to extract the meaning from pixels. Computer thinks like a machine
 - semantic info. → how many objects are
 - Eg. object categorizaⁿ, scene & context categorizaⁿ
 - vision as measurement device
 - stereo, structure from motion, reconstruction from photos.
- Object recogniⁿ → labelling each image with the dominant object.
 - + localizaⁿ with labelling
- Face detectⁿ + localizaⁿ
- Number plate reader
- Image segmentaⁿ → works like clustering of the pixels as per the objects.
- Object tracking in videos, activity recognition.
- Image captioning → image-to-text translation.

- depth from stereo \rightarrow recoder depth
by find coordinate x' to correspondingly x .
- depth from a single image, structure from motion (by rotating your camera)
- photometric stereo (camera fixed, shadow changes)
- Challenges
 1. viewpoint variation \rightarrow portrait
 - \hookrightarrow solve by augmentation.
 2. illumination \rightarrow too much Δ in lighting.
 3. scale \rightarrow bias issues
 4. deformation for animal poses.
 5. occlusion & background clutter
 6. motion blur & object inter-class variation.

L2 Deep learning for Computer Vision

- how does comp. recognise an image -
 - \rightarrow classifier
 - \rightarrow edges \rightarrow classifier
 - SIFT / HOG
 - \rightarrow edges \rightarrow histogram \rightarrow classifier



- we started deep learning features, it provides a richer solution space, thus, we can train it end-to-end learning by Back Prop.
 - Perceptron \rightarrow any funcⁿ that is linearly separable.
But what abt. XOR funcⁿ -
- x we can use more than 1 line. But 1 line can not separate.
 - x •
- then came, multi layer perceptrons.
 ↳ hidden layers, 'increase abstract'. hence, better to have more hidden layers than a single layer with large no. of neurons

L3 CNN (Convolutional Neural Networks)

- if you use a MLP to classify $3^3 224 \times 224$. typical a imagenet image \Rightarrow parameter explosion.
 ↳ it doesn't use to exploit local spatial information.

- CNN has more leverage —

1. local connectivity (LC)
2. parameter sharing (PS)
3. pooling / subsampling
4. ReLU (rectifier) non-linearity

- CNN are MLP with 2 constraints
→ LC & PS.

- MLNN (7×3)
21 params

MLNN - LC

($3 \times 3 = 9$ params)

hidden layer (3 nodes)

input layer (7 nodes)

m inputs, n output nodes & CNN

local connectivity of k nodes ($k < m$):

MLNN -

1. $m \times n$ params
2. $O(m \times n)$ runtime

MLNN - LC -

1. $k \times n$ params
2. $O(k \times n)$ runtime

2.3X runtime &
storage efficient.

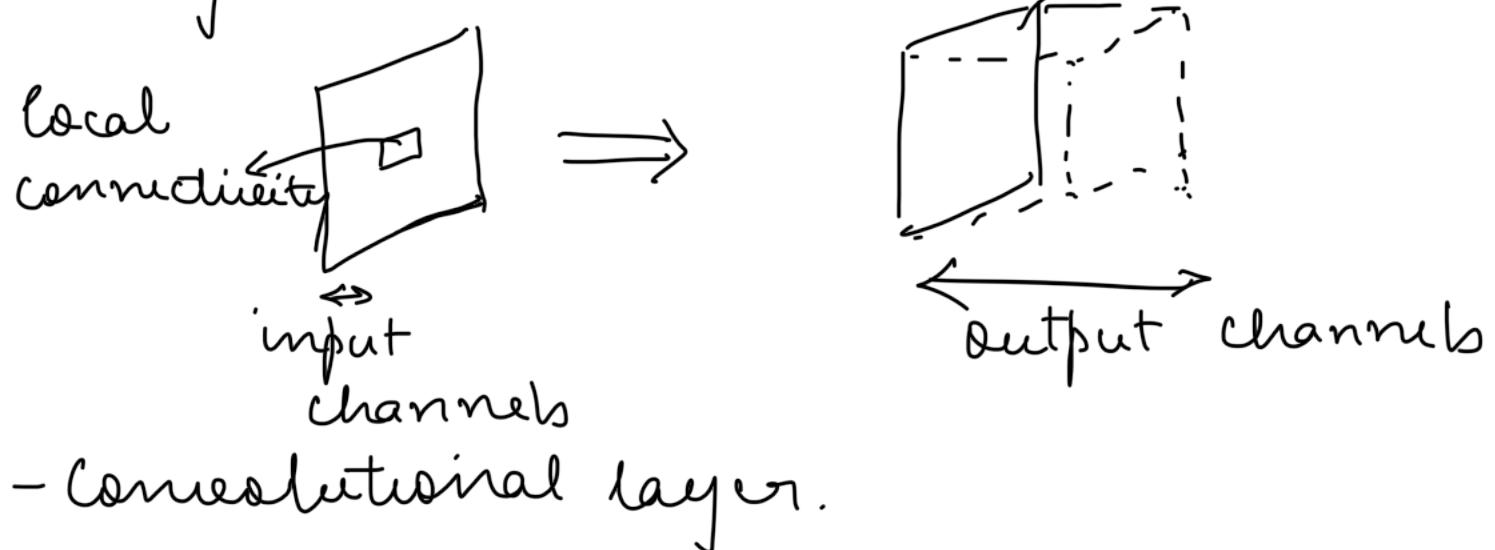
MLNN - LC - PS (3 params)

↳ 2.3X Faster & 7X storage
efficient.

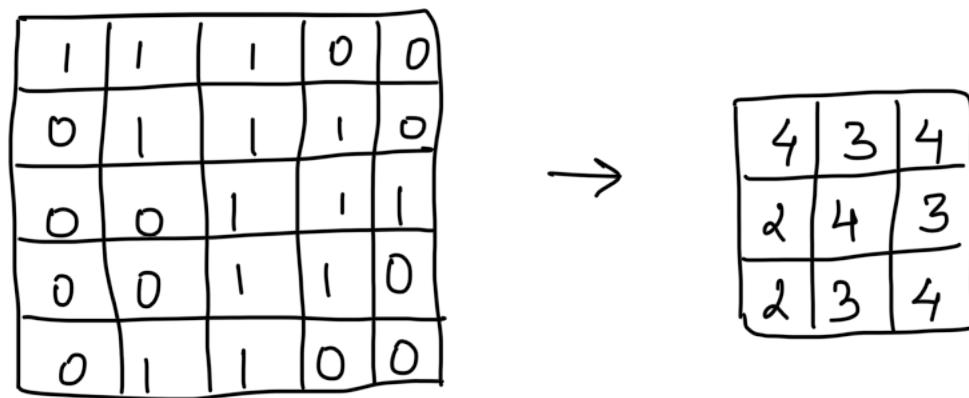
1. k params to store.

2. $O(k \times n)$ runtime.

- CNN with multiple input channels, output channels
- a generic CNN -



- Convolutional layer.



image

1. to reduce the no. of nets. (LC)
 2. to provide spatial invariance (PS)
- hyperparameters for CNN.
 1. 0 padding (to control input size spatially).
 2. stride (to produce smaller output volumes spatially)
 3. both padding & stride.
 4. single depth slice of max pool layer.
 - max pool layer -

1. to reduce the spatial size of the representation to reduce the amt. of param. and computa" in the network.
 2. avg. pooling or L2 pooling can be used, but not popular like max pooling.
- Activation func" → sigmoid func".
- ↳ It saturates & kills gradients (when the neuron is activated) at either tail of 0 or 1).
- ReLU (Rectified Linear Unit)
1. diminishes saturation & killing of gradients for the inputs.
 2. it accelerates convergence of SGD.
 3. tanh / sigmoid neurons involves expensive operations. (ReLU can be implemented by simply thresholding activation at 0.)
- Flattening → vectoriza" (converting $M \times N \times D$ tensor into $MND \times 1$)
- FC Layer → MLP, generally used in final layers, works as a classifier
- Softmax Layer → $Z_n = \frac{e^{x_n}}{\sum_{i=1}^k e^{x_i}}$

L4 Intro. to Object Recognition

- object recogⁿ (image classificaⁿ) - a core problem in CV.



→ classify it as an object.

- Challenges -

1. viewpoint variaⁿ → all the pixel value Δ s when the camera moves.
2. illuminaⁿ → too much contrast
3. deformaⁿ → diff. angles.
4. Occlusion & background clutter
5. interclass variation

- Now, we will study the CNN architectures

- LeNet 5

comes filters 5×5 , stride = 1

pooling layers were 2×2 , stride = 2

CONV - POOL - CONV - POOL - FC - FC.

L5 Alex Net & VGGNet

- had 8 layers

CONV,

MaxPool,

Norm1

CONV2

MaxPool2

Norm2

details

1. first use of ReLU
2. used Norm layers (now not in use)
3. heavy data augmentation

CONV3

4. dropout = 0.5

CONV4

5. batch size = 128

CONV5

6. SGD momentum = 0.9

MaxPool3

7. $\gamma = 0.01$

FC6

8. L_2 weight decay = $5e^{-6}$

FC7

9. 7 CNN ensemble.

FC8

- uses filters of 11×11 . It increases the no. of training parameters. Now, these days CNN uses cascade of filters of size 3×3 or 5×5 .

- VGGNet \rightarrow small filters, deep networks.

8 in Alex, 16 or 19 in VGG

11.7% error \Rightarrow 7.3%.

\rightarrow stack of 3 3×3 conv. (stride 1) layers has same effective receptive field as 1 7×7 conv layer. But now CNN becomes more deeper, more non-linearities. And lesser params to compute.

$3^2 \cdot 3^2 (c^2)$ vs $7^2 c^2$

if c is no. of channels per layer.

\rightarrow details -

1. no norm layers

2. uses ensembles for better result.

3. FC7 features generalizes well to other

tasks as well.

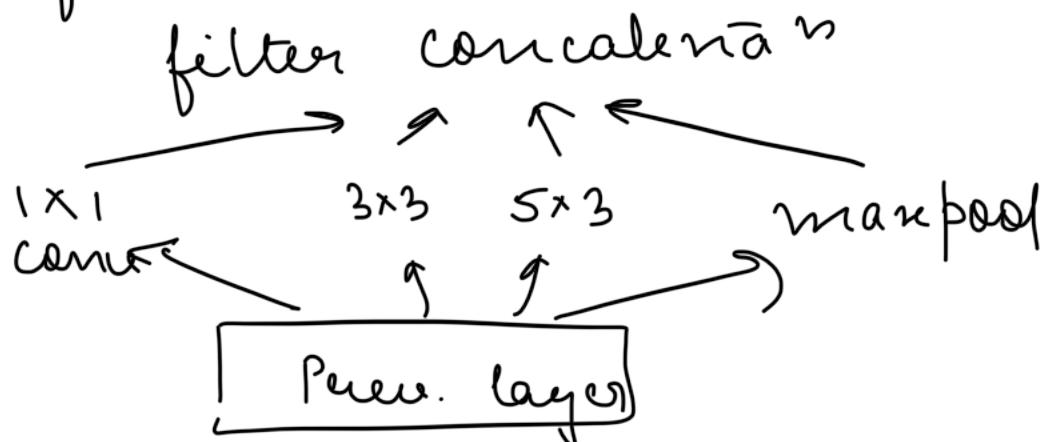
L6 Google Net & Res Net

- Google Net

→ Details → deeper networks with computational efficiency.

1. 22 layers.
2. efficient "Incep" module.
3. No FC layers.
4. only 5 M params. ($12 \times$ less than Alex)
5. error 6.7%.

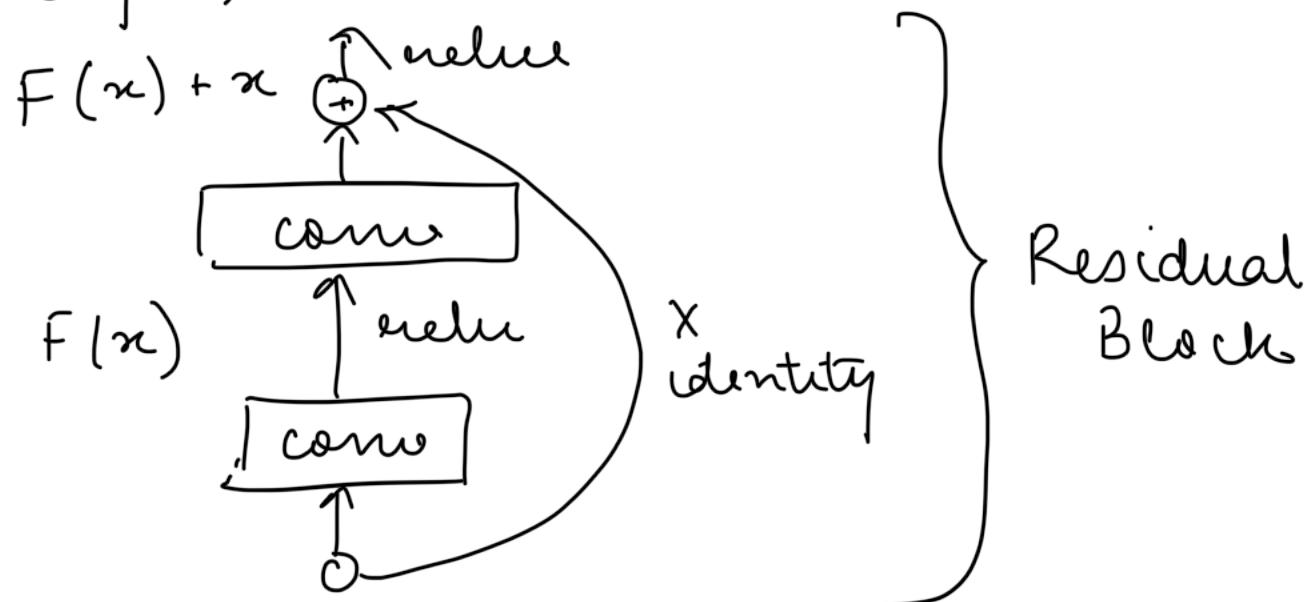
→ Incep" → design a good local network topology (network within a network) & then stack these modules on top of each other.



→ apply filter operations (11) on the inputs from prev. layers. Finally, concatenate all filter outputs together depth-wise.

→ problem → computa" complexity. There will be total 854 multiplica" operations

- ResNet - very deep networks using (152 layers) residual connections.



- the deeper model should be able to perform at least as well as the shallower model. A possible solution is by constructing the copied learning layers from the shallower model & setting additional layers to identity mapping.