

## Module 33: CNN architectures

→ what kind of tasks are CNNs used for? CNNs are generally applied to image related tasks.

(1) Classification: Predicting which class the image belongs to.

↓

Popular dataset: Imagenet (1000 categories)

↓

the bigger version (10000 categories)

↓

iconic photos: the class related object covers most space in the picture.

(2) Classification and Localization: Predicting the class, along with the location of object of the photo.

↓

The localization problem is a regression problem.

(3) Object detection: multiple objects in the photo.

The classification and localization for each object

(4) Instance segmentation: Contour of the object, rather than a box.

→ The decisions to be made for using CNN for image based tasks

i) No. of layers

ii) No. of filters in each layer

iii) filter size

iv) max pooling → when to use.



Not practical to try all permutations, rather the accepted practice is to use standard architectures which are known to give good results for a wide range of tasks.

→ The imagenet challenge: 1000 classes, with 1000 images each

shallow 2010 → top 5 error rate: 28.2%.

shallow 2011 → 25.8%.

8 2012 → Alexnet → 16.4% (drastic ↓)

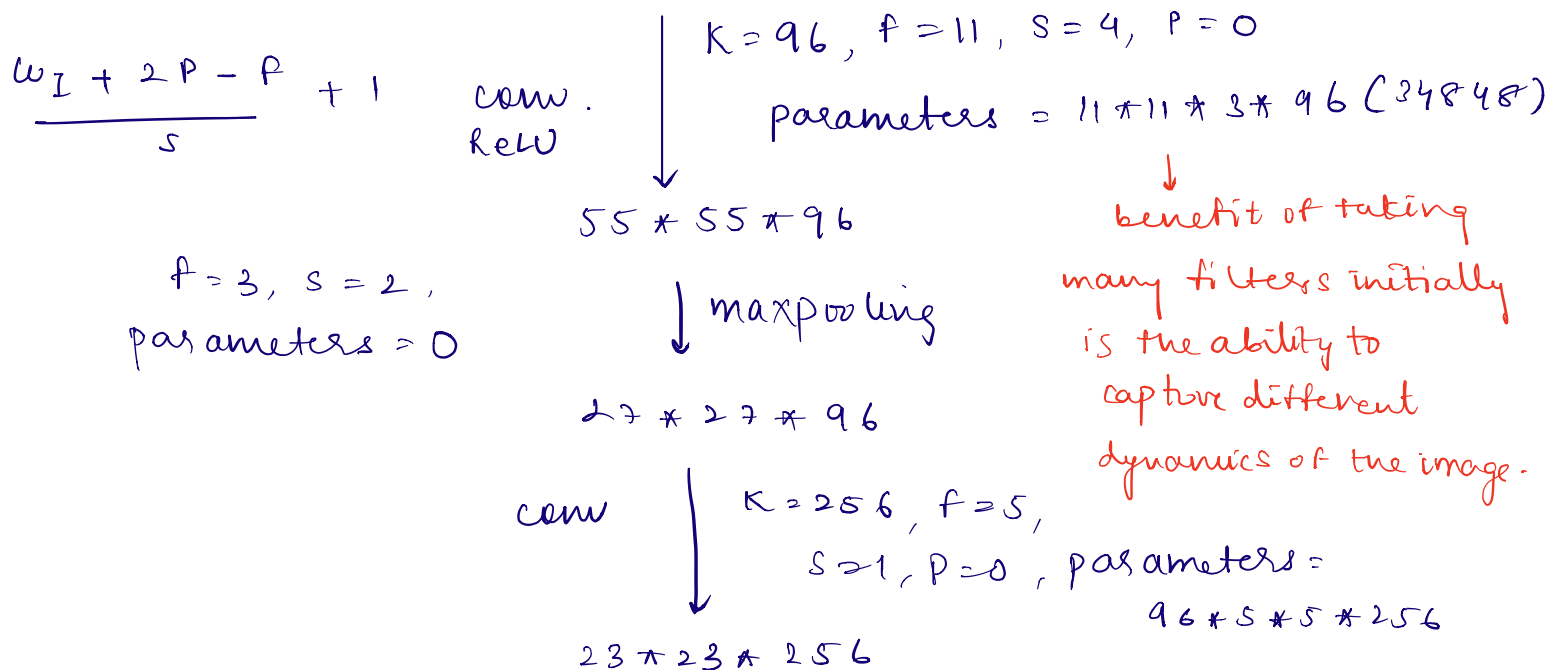
8 2013 → ZFNet → 11.7%.

19 2014 → VGG → 7.3%.

22 2014 → GoogLeNet → 6.7%.

152 2015 → Resnet → 3.57%.

→ Alexnet: inputs:  $227 * 227 * 3$  → channels



$$f=3, s=2$$

↓ maxpooling

$$11 \times 11 \times 256$$

maxpooling are not counted as layer, no parameters.

$$\text{conv} \downarrow \begin{array}{l} K=384, f=3, s=1, P=0, \\ \text{parameters} = 384 \times 3 \times 3 \times 256 \end{array}$$

$$9 \times 9 \times 384$$

$$\text{conv.} \downarrow \begin{array}{l} K=384, f=3, s=1, P=0 \\ \text{Parameters} = 384 \times 3 \times 3 \times 384 \end{array}$$

$$7 \times 7 \times 384$$

$$\text{conv} \downarrow \begin{array}{l} K=256, f=3, s=2, P=0 \\ \text{Parameters} = 384 \times 5 \times 3 \times 256 \end{array}$$

$$5 \times 5 \times 256$$

$$\text{maxpooling} \downarrow f=3, s=2$$

$$2 \times 2 \times 256$$

$$\text{flatten} \downarrow 4096$$

$$\text{dense} \downarrow 4096$$

$$\text{dense} \downarrow 1000 (\text{softmax})$$

<p>total parameters = 27.55m</p>
--------------------------------------

→ ZFnet : 8 layers again, can be understood using AlexNet, but differs at no. of filters, sizes, etc at some points

↓

(fewer parameters)

Layer 1:  $f = 11 \rightarrow 7$ , Difference in parameters = 20.7K  
 $w_0 = 56$ , scaled to 55 ( $55 \times 55 \times 96$ , same as AlexNet)

↓  
Layer 2 : No difference (Maxpooling)

↓  
Layer 3 : No difference

↓  
Layer 4 : No difference

↓  
Layer 5 :  $K = 384 \rightarrow 512$

Difference in parameters = 0.29 M  
(Increase in parameters)

↓  
Layer 6 :  $K = 384 \rightarrow 1024$

Difference in parameters = 0.8 M (+)

↓  
Layer 7 :  $K = 256 \rightarrow 512$

Difference in parameters = 0.36 M (+)

↓  
Layer 8 : No difference

↓  
Layer 9 : No difference

↓  
Layer 10 : No difference

↓  
Layer 11 : No difference

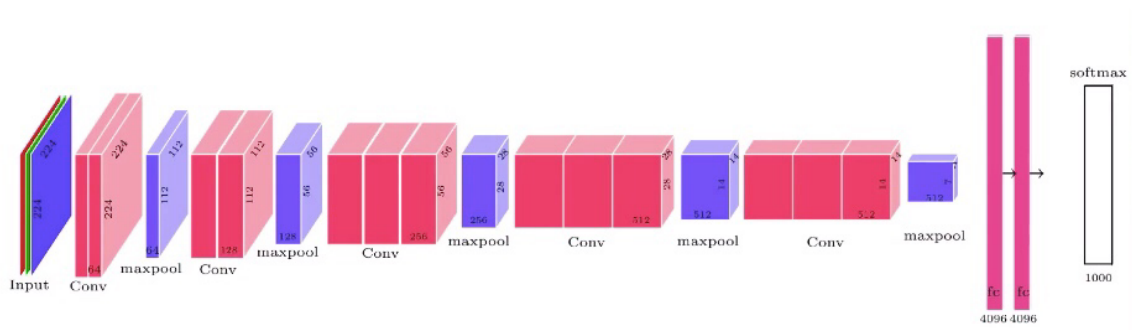
→ VGGNet (Replaced the notion of needing alternate conv and pooling layers, rather combining conv. layers and pooling layers)

↓  
Filters :  $3 \times 3$

↓

- 1) Back to back conv of :  $64 \times 3 \times 3$  (maintaining the dimensions along length and width)
- 2) maxpooling ( $\times 2$ )
- 3) Back to back conv. of :  $128 \times 3 \times 3$  ↗
- 4) maxpooling
- 5) 3 conv layers :  $256 \times 3 \times 3$
- 6) maxpooling
- 7) 3 conv layers :  $512 \times 3 \times 3$
- 8) maxpooling
- 9) 3 conv layers :  $512 \times 3 \times 3$
- 10) maxpooling
- 11) 2 fully connected layers (4096)
- 12) o/p layer (softmax)

VGG16



total parameters in non fc layers =  $\sim 16$  m

total parameters in fc layers =  $\sim 122$  m

most parameters are in the first fc layer ( $\sim 102$  m)