



3

## Agenda



11/30/2024

pra-sâmi

4

## Acknowledgement...

Geoffrey Everest Hinton CC FRS FRSC

- ❑ An English Canadian cognitive psychologist and computer scientist, most noted for his work on artificial neural networks.
- ❑ Since 2013, he divides his time working for Google (Google Brain) and the University of Toronto. In 2017, he cofounded and became the Chief Scientific Advisor of the Vector Institute in Toronto.
- ❑ With David Rumelhart and Ronald J. Williams, Hinton was co-author of a highly cited paper published in 1986 that popularized the **backpropagation algorithm** for training multi-layer neural networks, although they were not the first to propose the approach.
- ❑ Hinton is viewed as a **leading figure** in the deep learning community.
- ❑ The dramatic image-recognition milestone of the **AlexNet** designed in collaboration with his students Alex Krizhevsky and Ilya Sutskever for the ImageNet challenge 2012 was a breakthrough in the field of computer vision.

11/30/2024

pra-sâmi

5

What is Computer Vision...

11/30/2024

pra-sami

6

Ambulance given green light all through....

11/30/2024

pra-sami

7

Computer vision is making progress in leaps and bounds...

11/30/2024

pra-samí

8

Convolutional neural networks (CNN, ConvNet) is a class of deep, feed-forward (not recurrent) artificial neural networks that are applied to analyzing visual imagery

11/30/2024

pra-samí

9

## Computer Vision

- ❑ Self driving car
- ❑ Fully automated warehouse and ports
  - ❖ <https://youtu.be/BFV8ikY52iY>
- ❑ Image search services,
- ❑ Unlock phone
- ❑ Provide access to secure area
  - ❖ Open your house
  - ❖ Enter office without your access card
- ❑ Object identification Apps
  - ❖ Garment
  - ❖ Food,
  - ❖ Nature
- ❑ Natural style transfer
- ❑ Automatic video classification systems

11/30/2024

pra-sâmi

10

## Style Transfer



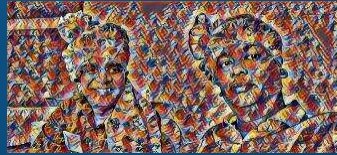
11/30/2024

pra-sâmi



11

## Style transfer



11/30/2024

pra-samī

12

## Style transfer



11/30/2024

pra-samī

13

## Computer Vision

- ❑ Have been used in image recognition since the 1980s
- ❑ Increase in computational power, the amount of available training data, CNNs have managed to achieve better performance
- ❑ Rapid advancement
  - ❖ Newer and Newer products and applications are coming up
  - ❖ Some of you will get a chance to directly work on these advance applications
- ❑ The development community is also very kind in sharing their success stories
- ❑ The ideas can be borrowed in other applications:
  - ❖ Voice recognition
  - ❖ Natural language processing (NLP)

11/30/2024

pra-sâmi

14

## Computer Vision

- ❑ What makes vision hard?
- ❑ Vision needs to be robust to a lot of transformations or distortions:
  - ❖ Change in pose/viewpoint
  - ❖ Change in illumination
  - ❖ Deformation
  - ❖ Occlusion (some objects are hidden behind others)
- ❑ Many object categories can vary wildly in appearance (e.g. chairs)

“Imaging a medical database in which the age of the patient sometimes hops to the input dimension which normally codes for weight!” - Geoff Hinton

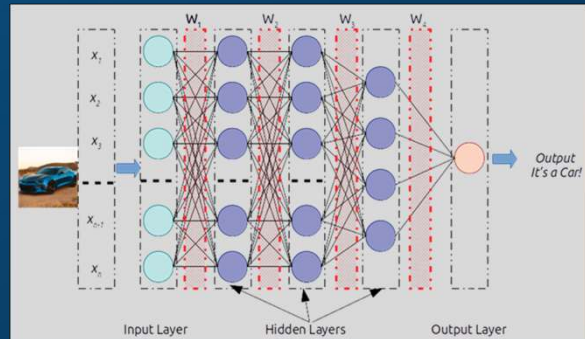
11/30/2024

pra-sâmi

15

## Why?

- ❑ Enough of sales pitch...
- ❑ Why not simply use a regular deep neural network with fully connected layers?
- ❑ Small (150 x 150 x 3) image has 67,500 pixels
- ❑ If we consider first hidden layer as 1000,
- ❑ First weight matrix ( $W_1$ ) will be 67,500 x 1000
- ❑ Do your math..... that size is huge



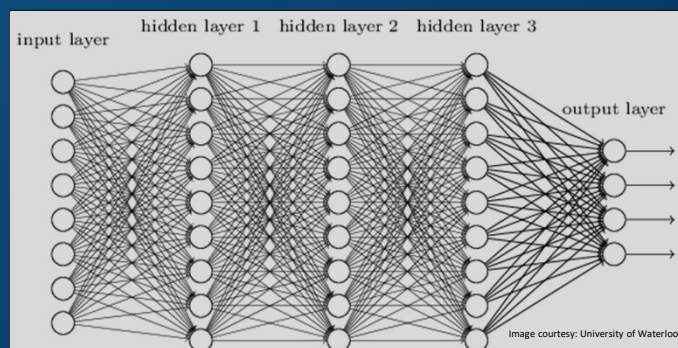
11/30/2024

pra-sami

16

## Smaller Network: CNN

- ❑ We know it is good to learn a small model
- ❑ Fully connected model, each hidden unit is processing every input
  - ❖ Do we really need all the edges?
- ❑ Can some of these be shared?



11/30/2024

pra-sami



17

Images are high-dimensional vectors. It would take a huge amount of parameters to characterize the network.

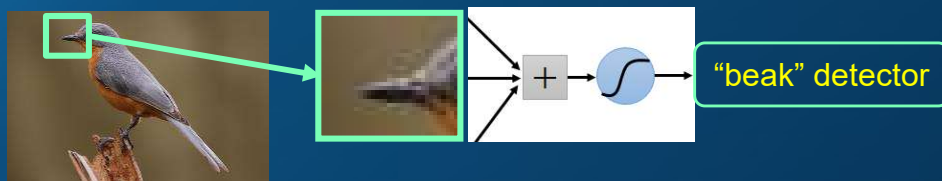
11/30/2024

pra-sâmi

18

## Learning an image...

- ❑ Some patterns are much smaller than the whole image
- ❑ Can represent a small region with fewer parameters



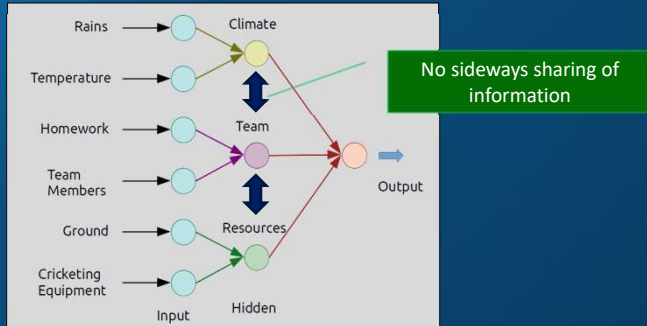
11/30/2024

pra-sâmi

19

## Learning an image...

- Same pattern appears in different places
  - ❖ Can they be compressed!
- What about training a lot of such “small” detectors and each detector must “move around”

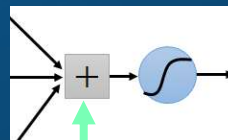
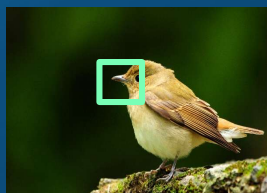
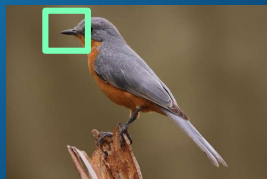


11/30/2024

pra-sâmi

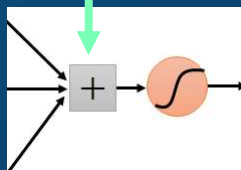
20

## Learning an image...



“upper-left beak”  
detector

They can be compressed to the same parameters.



“middle beak”  
detector

11/30/2024

pra-sâmi

21

## Learning an image...

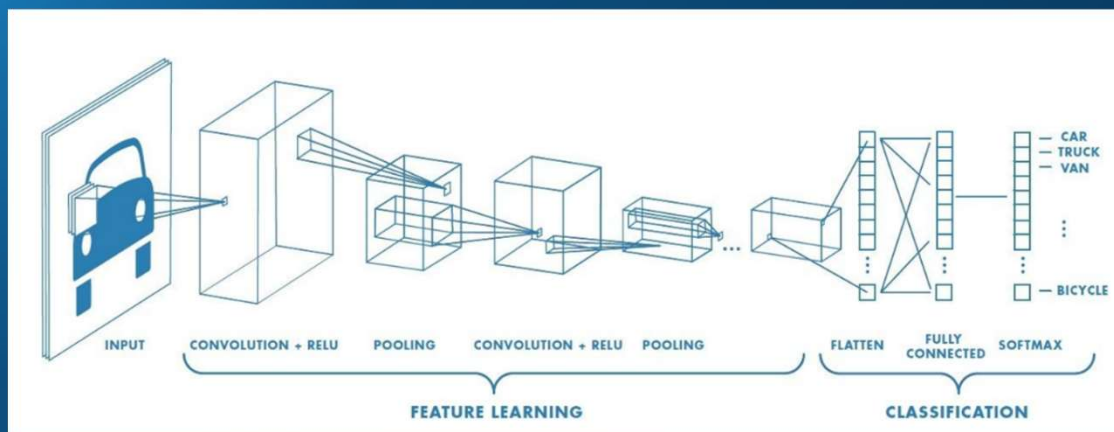
- ❑ The same sorts of features that are useful in analyzing one part of the image will probably be useful for analyzing other parts as well.
  - ❖ E.g., edges, corners, contours, object parts
- ❑ We want a neural net architecture that lets us learn a set of feature detectors that are applied at all image locations
- ❑ So far, we've seen a bunch of types of layers
  - ❖ Fully connected layers (dense)
  - ❖ Embedding layers (i.e. lookup tables)
  - ❖ A few more in RNNs (GRU, LSTMs, etc.)
- ❑ Different layers could be stacked together to build powerful models
- ❑ Let's add another set of layers: the convolution layer, pooling layer...

11/30/2024

pra-sami

22

## Overall Layout



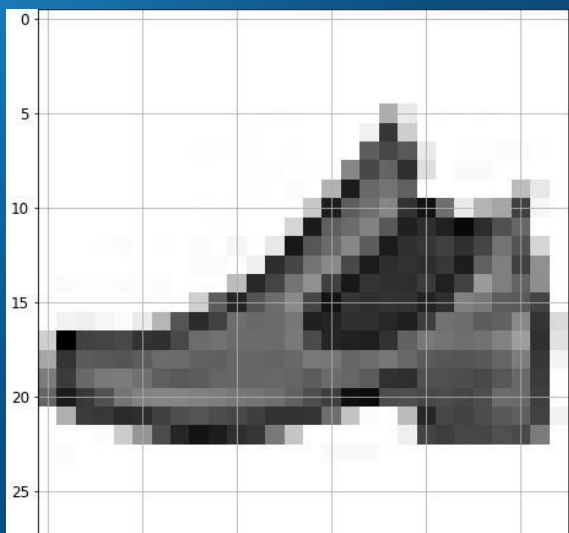
11/30/2024

pra-sami

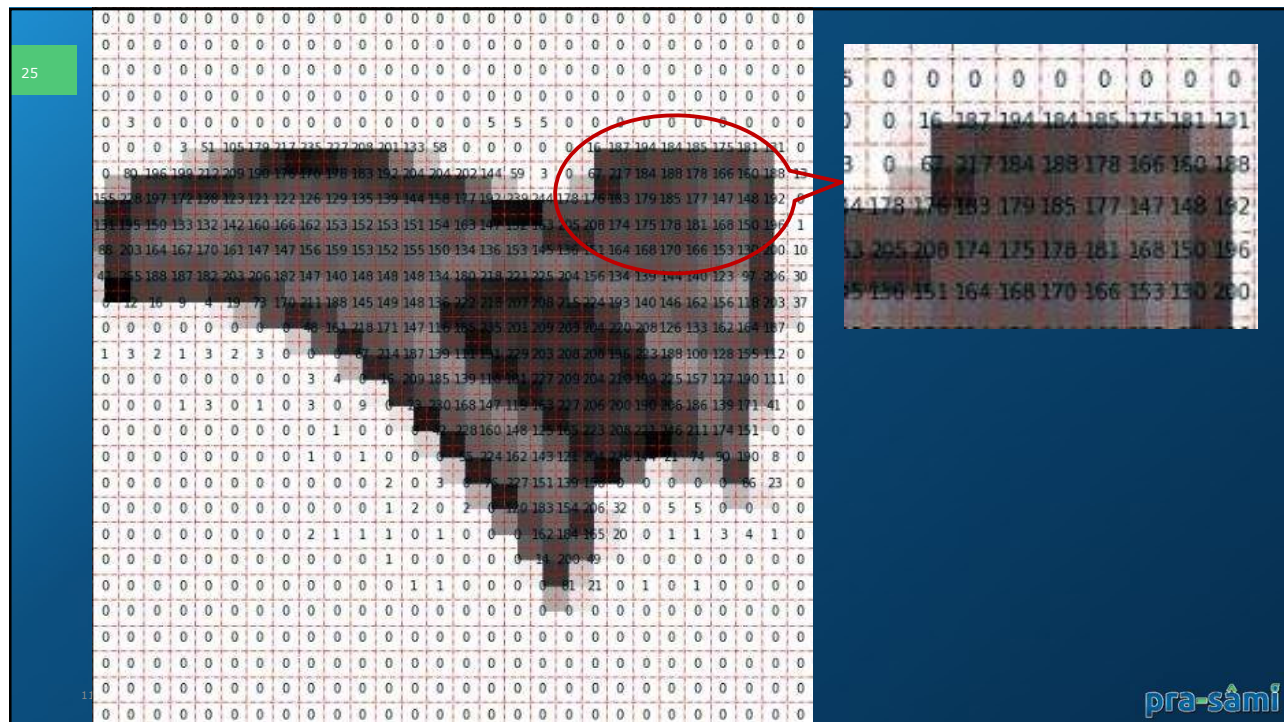
Each column of hidden units looks at a small region of the image, and the weights are shared between all image locations.

pra-sami

## What computer sees...

[illegible]

pra-sâmi



26

## A Convolutional Layer

- ❑ A CNN is a neural network with some convolutional layers
  - ❖ And, of course, a few other layers
- ❑ A convolutional layer has a number of filters that does convolutional operation
  - ❖ Some of the literature would call it Kernel

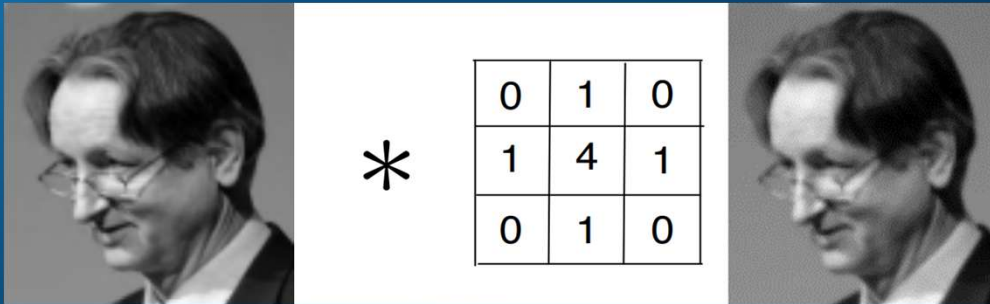
11/30/2024

pra-sami



27

What does this Convolution Filter/ Kernel do?

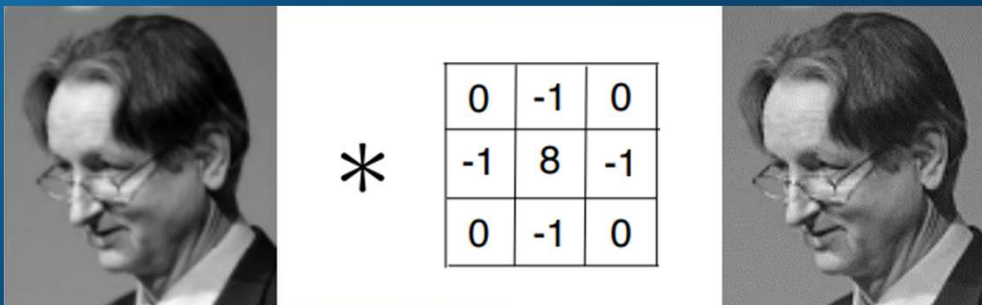


11/30/2024

pra-samí

28

What does this Convolution Filter/ Kernel do?

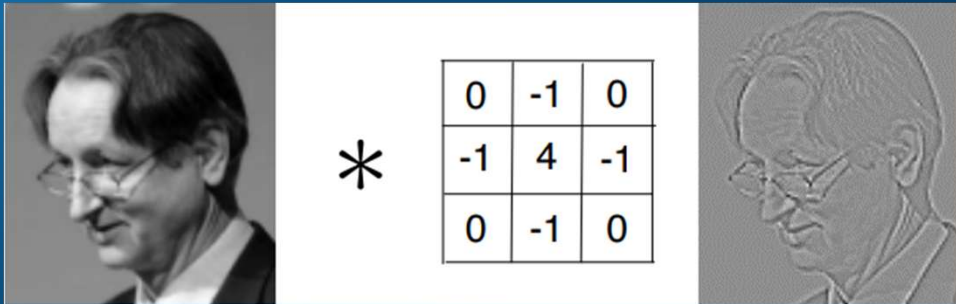


11/30/2024

pra-samí

29

What does this Convolution Filter/ Kernel do?

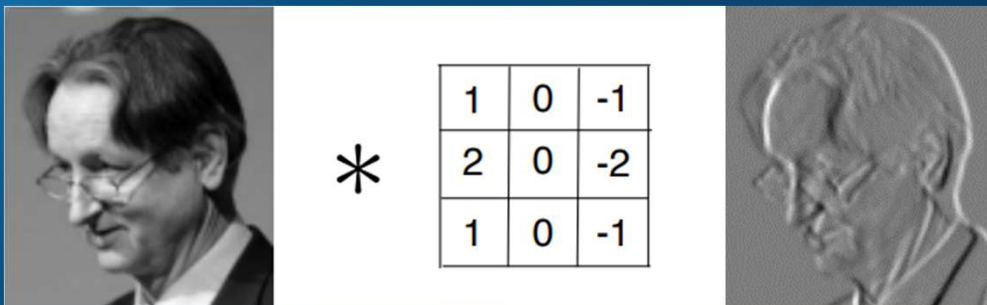


11/30/2024

pra-sami

30

What does this Convolution Filter/ Kernel do?



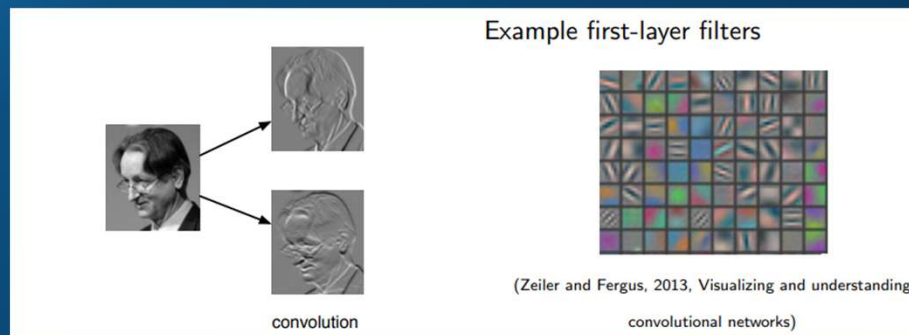
11/30/2024

pra-sami

31

## Convolutional Networks

- Two kinds of layers:
  - ❖ Detection layers (or convolution layers)
  - ❖ Pooling layers
- The convolution layer has a set of filters.
  - ❖ Output is a set of feature maps, each one obtained by convolving the image with a filter.



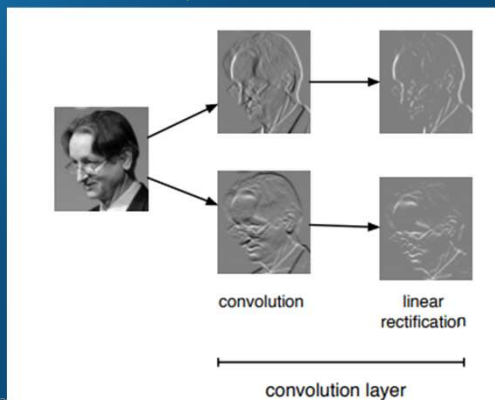
11/30/2024

pra-samī

32

## Convolutional Networks

- It's common to apply a linear rectification (activations) nonlinearity or even something else:
  - ❖  $y_i = \text{Relu}(z_i)$ ,
  - ❖ May be,  $\text{Tanh}(z_i)$ , etc.



11/3

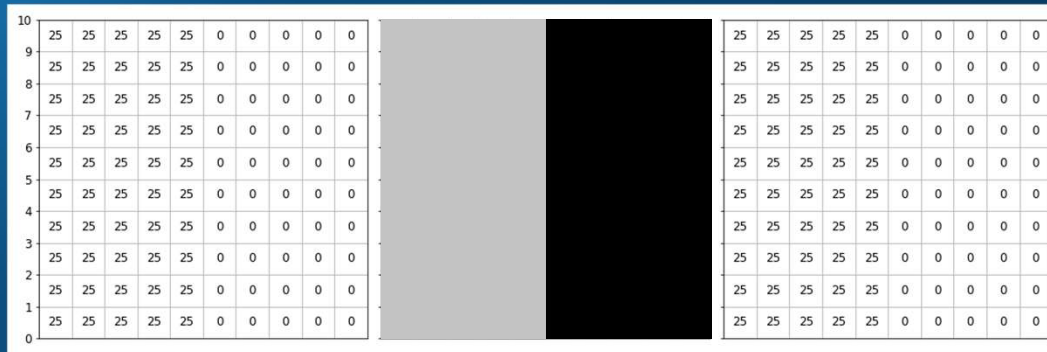
- Convolution is a linear operation
- Therefore, we need a nonlinearity:
  - ❖ Otherwise two convolution layers would be no more powerful than one
- Two edges in opposite directions shouldn't cancel
- Non-linearity makes the gradients sparse, which helps optimization

pra-samī

33

## Image with a edge

- Convolution is basic building block of image recognition
- Using edge detection as an example in following image...



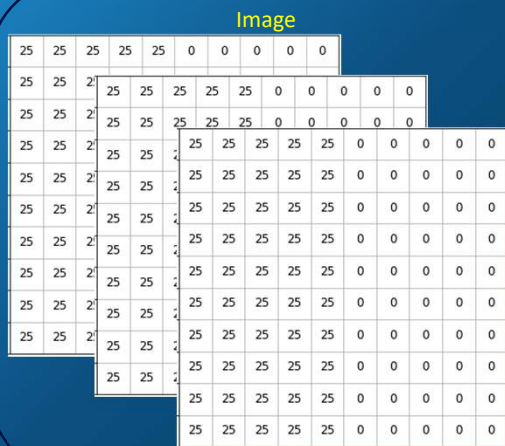
- Apply filters on the image!

11/30/2024

pra-sami

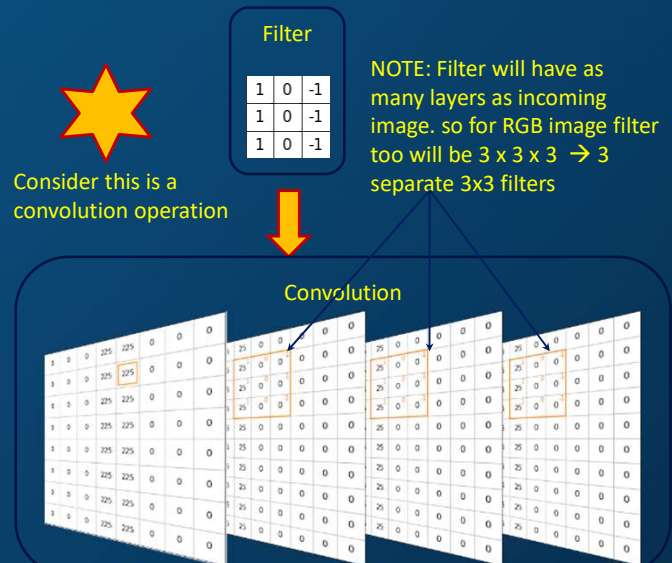
34

## Convolution on 3D images (RGB )



If you are looking for edge in one channel only, make rest of them as zeros.

11/30/2024



pra-sami

35

## Convolution on 3D images (RGB )

### □ First convolution

25	<sup>1</sup>	25	<sup>0</sup>	25	<sup>-1</sup>	25	25	0	0	0	0	0
25	<sup>1</sup>	25	<sup>0</sup>	25	<sup>-1</sup>	25	25	0	0	0	0	0
25	<sup>1</sup>	25	<sup>0</sup>	25	<sup>-1</sup>	25	25	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0	0	0

### □ Layer R

$$\diamond = 25*1 + 25*1 + 25*1 + 0 + 0 + 0 - 1*25 - 1*25 - 1*25$$

$$\diamond = 0$$

### □ Layer G

$$\diamond = 25*1 + 25*1 + 25*1 + 0 + 0 + 0 - 1*25 - 1*25 - 1*25$$

$$\diamond = 0$$

### □ Layer B

$$\diamond = 25*1 + 25*1 + 25*1 + 0 + 0 + 0 - 1*25 - 1*25 - 1*25$$

$$\diamond = 0$$

$$\square \text{ Total} = 0 + 0 + 0 = 0$$

11/30/2024

pra-sâmi

36

## Convolution on 3D images (RGB )

### □ Second convolution

✧ It will be identical to First

25	25	<sup>1</sup>	25	<sup>0</sup>	25	<sup>-1</sup>	25	0	0	0	0	0
25	25	<sup>1</sup>	25	<sup>0</sup>	25	<sup>-1</sup>	25	0	0	0	0	0
25	25	<sup>1</sup>	25	<sup>0</sup>	25	<sup>-1</sup>	25	0	0	0	0	0
25	25	25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	25	25	0	0	0	0	0	0

### □ Layer R

$$\diamond = 25 + 25 + 25 + 0 + 0 + 0 - 25 - 25 - 25$$

$$\diamond = 0$$

### □ Layer G

$$\diamond = 25 + 25 + 25 + 0 + 0 + 0 - 25 - 25 - 25$$

$$\diamond = 0$$

### □ Layer B

$$\diamond = 25 + 25 + 25 + 0 + 0 + 0 - 25 - 25 - 25$$

$$\diamond = 0$$

$$\square \text{ Total} = 0 + 0 + 0 = 0$$

11/30/2024

pra-sâmi



37

## Convolution on 3D images (RGB )

□ What happens 4<sup>th</sup> step

25	25	25	25 <sup>1</sup>	25 <sup>0</sup>	0 <sup>-1</sup>	0	0	0	0
25	25	25	25 <sup>1</sup>	25 <sup>0</sup>	0 <sup>-1</sup>	0	0	0	0
25	25	25	25 <sup>1</sup>	25 <sup>0</sup>	0 <sup>-1</sup>	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0

□ Layer R

$$\diamond = 25 + 25 + 25 + 0 + 0 + 0 - 0 - 0 - 0$$

$$\diamond = 75$$

□ Layer G

$$\diamond = 25 + 25 + 25 + 0 + 0 + 0 - 0 - 0 - 0$$

$$\diamond = 75$$

□ Layer B

$$\diamond = 25 + 25 + 25 + 0 + 0 + 0 - 0 - 0 - 0$$

$$\diamond = 75$$

□ Total = 75 + 75 + 75 = 225

11/30/2024

pra-sami

38

## Convolution on 3D images (RGB )

□ And for 5<sup>th</sup> Step

25	25	25	25	25 <sup>1</sup>	0 <sup>0</sup>	0 <sup>-1</sup>	0	0	0
25	25	25	25	25 <sup>1</sup>	0 <sup>0</sup>	0 <sup>-1</sup>	0	0	0
25	25	25	25	25 <sup>1</sup>	0 <sup>0</sup>	0 <sup>-1</sup>	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0

□ Layer R

$$\diamond = 25 + 25 + 25 + 0 + 0 + 0 - 0 - 0 - 0$$

$$\diamond = 75$$

□ Layer G

$$\diamond = 25 + 25 + 25 + 0 + 0 + 0 - 0 - 0 - 0$$

$$\diamond = 75$$

□ Layer B

$$\diamond = 25 + 25 + 25 + 0 + 0 + 0 - 0 - 0 - 0$$

$$\diamond = 75$$

□ Total = 75 + 75 + 75 = 225

11/30/2024

pra-sami

39

## Convolution on 3D images (RGB )

□ 6<sup>th</sup> Step onwards again all values are 0

25	25	25	25	25	0 <sup>1</sup>	0 <sup>0</sup>	0 <sup>-1</sup>	0	0
25	25	25	25	25	0 <sup>1</sup>	0 <sup>0</sup>	0 <sup>-1</sup>	0	0
25	25	25	25	25	0 <sup>1</sup>	0 <sup>0</sup>	0 <sup>-1</sup>	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0

□ Layer R

$$\diamond = 0 + 0 + 0 + 0 + 0 + 0 - 0 - 0 - 0$$

$$\diamond = 0$$

□ Layer G

$$\diamond = 0 + 0 + 0 + 0 + 0 + 0 - 0 - 0 - 0$$

$$\diamond = 0$$

□ Layer B

$$\diamond = 0 + 0 + 0 + 0 + 0 + 0 - 0 - 0 - 0$$

$$\diamond = 0$$

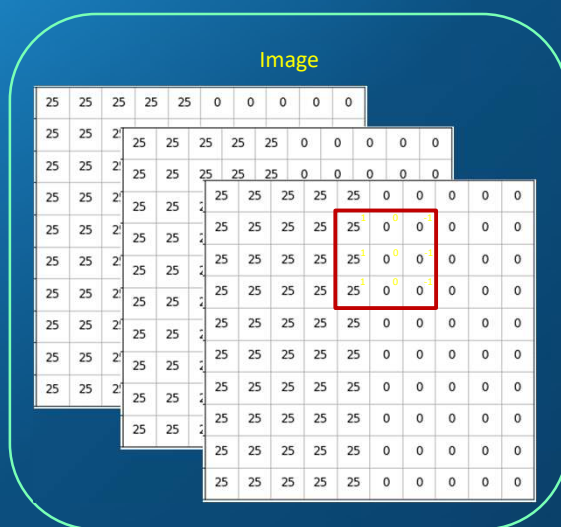
□ Total = 0 + 0 + 0 = 0

11/30/2024

pra-sami

40

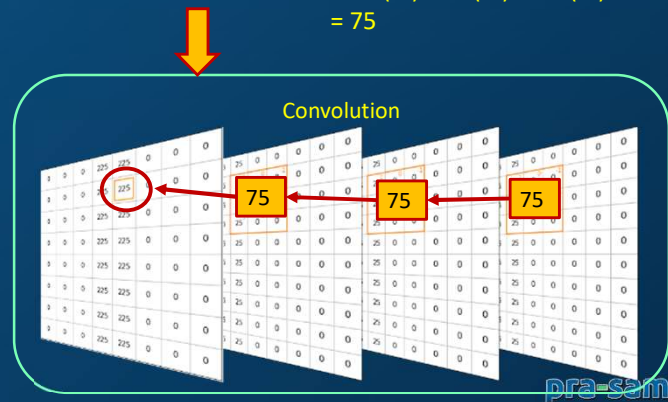
## Convolution



Filter

1	0	-1
1	0	-1
1	0	-1

$$\begin{aligned} \square & 25 * 1 + 25 * 1 + 25 * 1 \\ & + 0 * 0 + 0 * 0 + 0 * 0 \\ & + 0 * (-1) + 0 * (-1) + 0 * (-1) \\ & = 75 \end{aligned}$$



11/30/2024

pra-sami

41

## Convolution

Image

25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0

11/30/2024

Filter

1	0	-1
1	0	-1
1	0	-1

- Every filter we deploy results in 2D matrix

Convolution

pra-sami

42

## Convolution

- How many steps filter can take before it goes out of image?
- $10 - 3 + 1 = 8$  in either direction...

25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0

11/30/2024

pra-sami

43

## Convolution

$$\square 10 - 3 + 1 = 8 \times 8$$

25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0

□ Output will be a 2D Matrix

□ Assuming it moves by one step

□ Given that size of the image is 10 and size of the filter is 3

□ Taking :  $10 - 3 + 1 = 8$  steps

□ Output image will have 8 cells

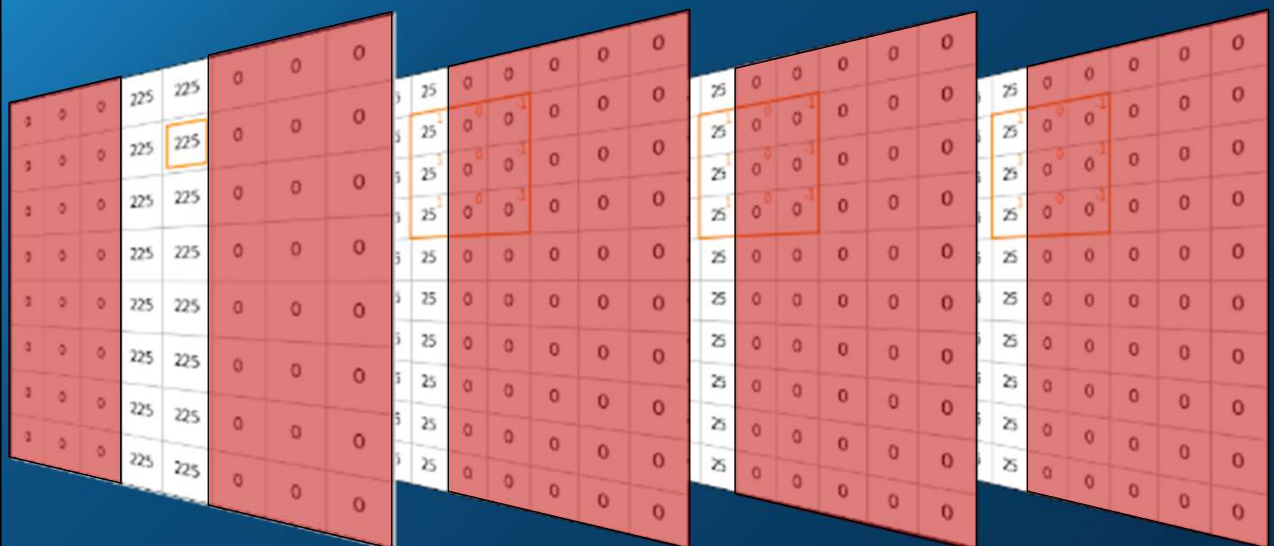
□ Hence, Output image size  
 $= \{ nH_{in} - nF + 1 \} \times \{ nW_{in} - nF + 1 \}$

11/30/2024

pra-sami

44

## Has it detected the edge?



11/30/2024

pra-sami

45

What if we move two steps at a time

11/30/2024

pra-sâmi

46

## Stride

25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0

- ❑ Steps are called stride
- ❑ If we move 2 steps at a time, we can move 4 steps only
- ❑ In other word with stride of 2
  - ❖ Given that size of the image is 10 and size of the filter is 3
- ❑ Output image size will be :  $(10 - 3)/2 + 1 = 4.5$  or 4
- ❑ For fractions pick lower integer
- ❑ Over flow not permitted

11/30/2024

pra-sâmi



47

## Stride

$$\square (10 - 3)/2 + 1 = 4 \times 4$$

25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0

□ Steps are called stride

$$\square \text{ Hence, Output image size} = \frac{\{ (nH_{in} - nF) / \text{stride} + 1 \} \times \{ (nW_{in} - nF) / \text{stride} + 1 \}}$$

- If we apply multiple filters → this layer will have 3D matrix.
- ❖ Each layer corresponding to one filter.

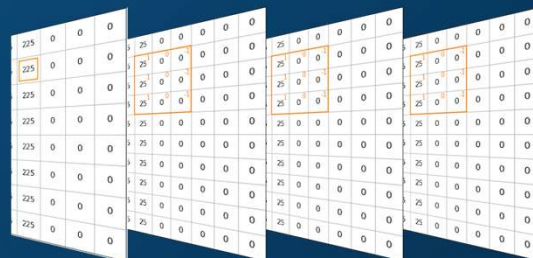
11/30/2024

pra-sami

48

## Convolution

- Apply filters and stack filtered layers together to make a 3D matrix
- Hence from 3 layer RGB, we can construct as many layers as number of filters applied...
- Move "stride" steps, generally one or two
  - ❖ one in most cases...
- Strongly advisable to keep filters as odd shape (3,3) or (5,5)
- Two strong reason..
- We do not want asymmetric padding
  - ❖ Not good for learning features
  - ❖ It's better to have central point of the filter



11/30/2024

pra-sami



51

## Another Convolution

- Layer R =  $3 + 3 + 102 - 98 - 206 - 252 = -448$
- Layer G =  $51 + 51 + 169 - 163 - 196 - 3 = -91$
- Layer B =  $252 + 252 + 203 - 207 - 10 - 18 = 472$
- Total =  $-448 + -91 + 472 = -67$

The diagram illustrates a convolution operation across three layers (R, G, B) and a final output grid. The layers are represented as 3D planes with numbers, and the output is a 2D grid.

3	3	98	140	51	51	163	225	252	252	207	131	-67	23
3	170	206	206	51	178	196	196	252	147	10	10	313	343
102	220	252	252	169	170	3	3	203	13	18	18	559	390
145	220	252	252	233	170	3	3	113	13	18	18	498	390
158	220	252	252	252	170	3	3	3	13	18	18	498	390
158	220	252	252	252	170	3	3	3	13	18	18	498	390
145	220	252	252	233	170	3	3	113	13	18	18	559	390
102	220	252	252	169	170	3	3	203	13	18	18	318	344
3	177	213	213	51	172	182	182	252	142	12	12	-62	24
3	3	98	140	51	51	163	225	252	252	207	131		

11/30/2024

pra-sami

52

## Another Convolution

- Incoming Image shape = (10,10,3)
  - ✦  $nH_{in} = 10$  ;  $nW_{in} = 10$ ,  $nC = 3$
- Filter shape = ( 3, 3, 3)
  - ✦  $nF = 3$  ;  $nF = 3$ ,  $nC = 3$
- Stride = 1
- Hence the size will be 8 x 8 after convolution

The diagram illustrates a convolution operation across three layers (R, G, B) and a final output grid. The layers are represented as 3D planes with numbers, and the output is a 2D grid.

3	3	98	140	51	51	163	225	252	252	207	131	-67	23
3	170	206	206	51	178	196	196	252	147	10	10	313	343
102	220	252	252	169	170	3	3	203	13	18	18	559	390
145	220	252	252	233	170	3	3	113	13	18	18	498	390
158	220	252	252	252	170	3	3	3	13	18	18	498	390
158	220	252	252	252	170	3	3	3	13	18	18	498	390
145	220	252	252	233	170	3	3	113	13	18	18	559	390
102	220	252	252	169	170	3	3	203	13	18	18	318	344
3	177	213	213	51	172	182	182	252	142	12	12	-62	24
3	3	98	140	51	51	163	225	252	252	207	131		

11/30/2024

pra-sami

53

## Another Convolution

- Single filter convolution:
- Layer R =  $3 + 3 + 102 - 98 - 206 - 252 = -448$
- Layer G =  $51 + 51 + 169 - 163 - 196 - 3 = -91$
- Layer B =  $252 + 252 + 203 - 207 - 10 - 18 = 472$
- Total =  $-448 + -91 + 472 = -67$

- Incoming Image shape = (10,10,3)

❖ nHin = 10 ; nWin = 10, nC = 3

- Filter shape = ( 3, 3, 3)

❖ → nF = 3 ; nF = 3, nC = 3

- Stride = 1

- Hence the size will be 8 x 8 after convolution

In convolution,:

- With every convolution image is shrinking
- Corners and edges of image are used less frequently than the middle

11/30/2024

pra-sâmi

54

## Other filters

- We have seen vertical filter... How about horizontal Filter....

- No surprises there....

1	1	1
0	0	0
-1	-1	-1

- The math will be exactly the same and we would get horizontal edge

11/30/2024

pra-sâmi

55

## Horizontal Edge...



11/30/2024

pra-sami

56

## Other filters

□ Sobel Filter...

1	0	-1
2	0	-2
1	0	-1

□ There was a lot of debate on filters....

□ Researchers kept trying various numbers...

□ Why not learn these parameters...

$w_1$	$w_4$	$w_7$
$w_2$	$w_5$	$w_8$
$w_3$	$w_6$	$w_9$

11/30/2024

pra-sami



57

## Shade Reversal

- So far we have seen lighter to darker shade filters...
- What happens if we move from darker shade to lighter

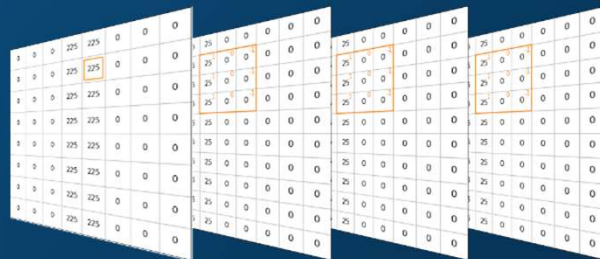
11/30/2024

pra-sami

58

## Shade Reversal

- So far we have seen lighter to darker shade filters...
- What happens if we move from darker shade to lighter
- We will again get the edge only it will be negative this time...



11/30/2024

pra-sami

59

## Convolving a Volume

- So far we have shown that same filter is applied to all layers
- In theory, it is possible to have a filter which is looking for edges in red channel alone...

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

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

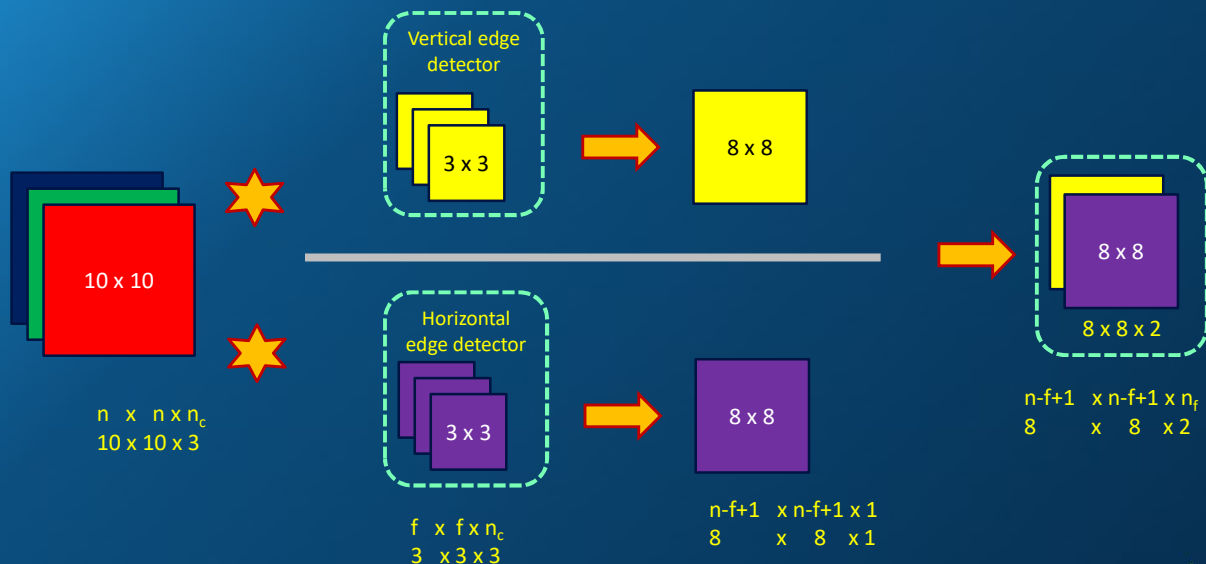
- So far we have been showing that 3D image converts to 2D image when we apply filter
- By applying a number of filters to detect different edges, we can have 3d Convolutional Volumes.

11/30/2024

pra-sâmi

60

## Multiple Filters



11/30/2024

pra-sâmi

61

## Two Issues with the convolution...

- ❑ With every convolution image is shrinking
  - ❖ Knowing that 100s of layer is not uncommon in the architecture
  - ❖ Image can soon become 1px X 1px
- ❑ Corners and edges of image are used less frequently than the middle

11/30/2024

pra-sâmi

62

What if we zero pad the image all around... will it help?

11/30/2024

pra-sâmi



65

## Convolution after Padding

- Incoming image shape = ( 10, 10, 3 )
  - ❖ i.e.  $nH_{in} = 10$  ;  $nW_{in} = 10$  ;  $nC = 3$
- Padding  $p = 1$
- Padded image shape = ( 12, 12, 3 )
  - ❖ i.e.  $nH_{in} = 12$  ;  $nW_{in} = 12$  ;  $nC = 3$
- Filter shape = ( 3, 3, 3 )
  - ❖ i.e.  $nF = 3$  ;  $nF = 3$ ,  $nC = 3$
- Assuming we move "stride" steps at any time
  - ❖ i.e. stride = 1

- Output image size:

$$= \left\{ \frac{nH_{in} - nF + 2 * p}{stride} + 1 \right\}$$

x

$$\left\{ \frac{nW_{in} - nF + 2 * p}{stride} + 1 \right\}$$

$$\begin{aligned} \square \text{ Image Size} &= \left\{ \frac{10 - 3 + 2 * 1}{1} + 1 \right\} \\ &\quad \times \\ &\quad \left\{ \frac{10 - 3 + 2 * 1}{1} + 1 \right\} \\ &= 10 \times 10 \end{aligned}$$

We are back to original size...

11/30/2024

pra-sâmi

66

## How much to pad???

- There are two recommended mechanism
- **Valid** : output is calculated as
 
$$\left\{ \frac{nH_{in} - nF + 2 * p}{stride} + 1 \right\} \times \left\{ \frac{nW_{in} - nF + 2 * p}{stride} + 1 \right\}$$
- So for 10 x 10 image a 5 x 5 filter with 1 px padding, image size will be 8 x 8
- **Same** : do the padding in such a way so that resultant image is of same size
 
$$\left\{ \frac{nH_{in} - nF + 2 * p}{stride} + 1 \right\} \times \left\{ \frac{nW_{in} - nF + 2 * p}{stride} + 1 \right\} = nH_{in} \times nW_{in}$$
  - ❖ or  $p = (nF - 1) / 2$  for stride = 1

11/30/2024

pra-sâmi

67

## How much to pad???

- With  $p = (nF - 1)/2$  for  $\text{stride} = 1$ ;
- We want  $p$  to be an integer and hence
  - ❖ Need  $nF$  to be odd
- For even value of  $nF$  we would end up in asymmetric padding.
- Unless we feel one edge of the image is more important than other, there is no need to have asymmetric padding

11/30/2024

pra-sâmi

68

## Cross-Correlation vs. Convolution

11/30/2024

pra-sâmi



69

## Cross-Correlation vs. Convolution

- ❑ In Signal Theory and Maths
- ❑ Convolution involves multiplying the filter after mirroring on both axis
- ❑ i.e. for filter  $\begin{bmatrix} 3 & 4 & 5 \\ 1 & 0 & 2 \\ -1 & 9 & 7 \end{bmatrix}$
- ❑ It will be mirrored along both axis...  $\begin{bmatrix} 7 & 9 & -1 \\ 2 & 0 & 1 \\ 5 & 4 & 3 \end{bmatrix}$
- ❑ Then we do element wise multiplication.
- ❑ Signal Engineers will agree with me... ☺
- ❑ Such correlations have properties like associative  $(a*b)*c = a*(b*c)$  and all other properties

11/30/2024

pra-sâmi

70

## Cross-Correlation vs. Convolution

- ❑ So that we are correct semantically...
- ❑ What we are doing is called Cross-Correlation....
- ❑ However, Data Scientists across the world have been using filters without reversing it and still call it Convolution...

Now you know... don't write home about it... ☺

11/30/2024

pra-sâmi

71

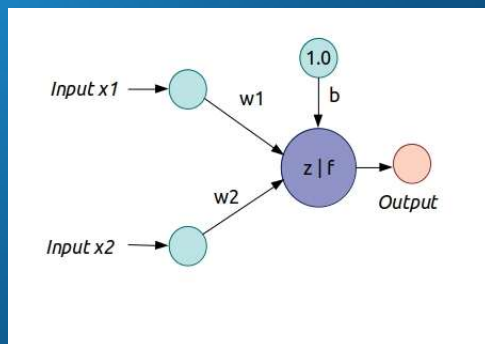
## One layer of Convolutional Net

11/30/2024

pra-sâmi

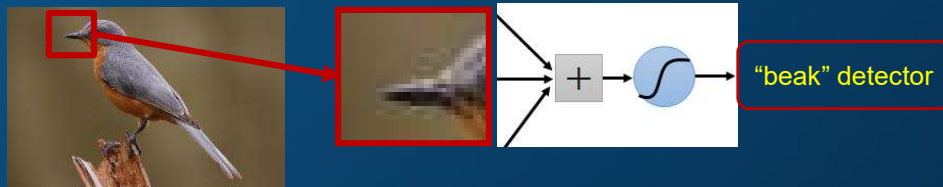
72

## One layer of Conv Net



$$Z_1 = a_0 \cdot W_1 + b_1$$

$$a_1 = \text{Relu}(Z_1)$$

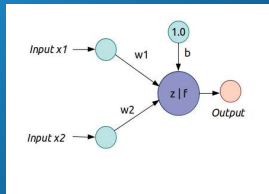


11/30/2024

pra-sâmi

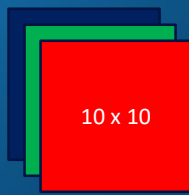
73

## One Layer of Conv Net



$$Z_1 = a_0 \cdot W_1 + b_1$$

$$a_1 = \text{Relu}(Z_1)$$

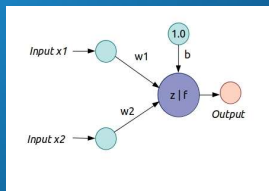


11/30/2024

pra-sâmi

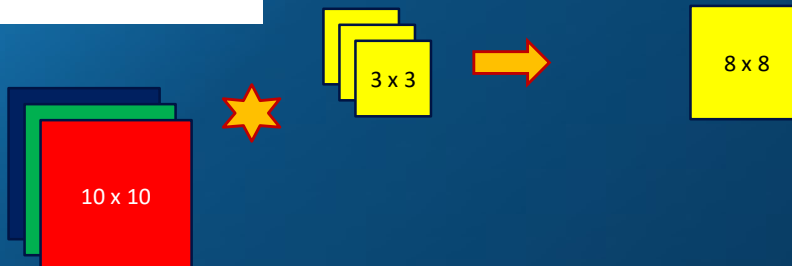
74

## One Layer of Conv Net



$$Z_1 = a_0 \cdot W_1 + b_1$$

$$a_1 = \text{Relu}(Z_1)$$

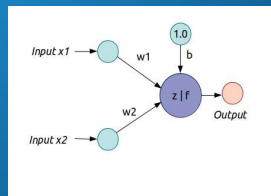


11/30/2024

pra-sâmi

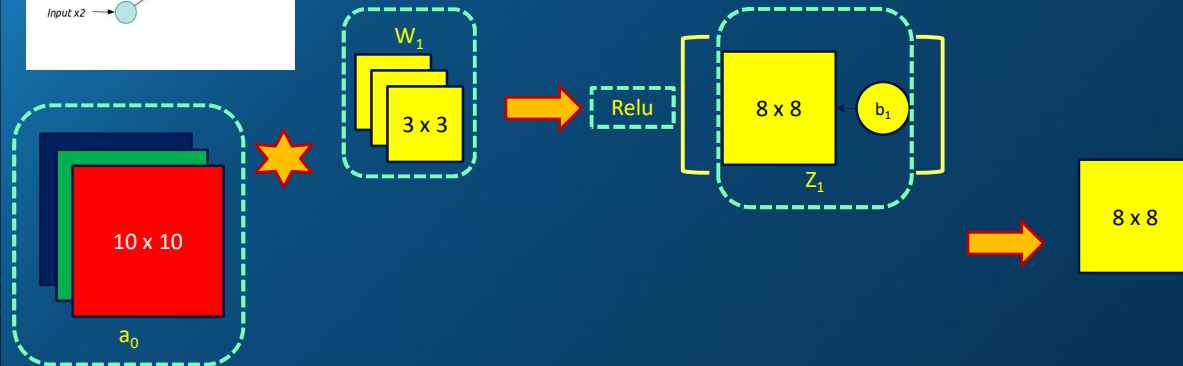
75

## One Layer of Conv Net



$$Z_1 = a_0 \cdot W_1 + b_1$$

$$a_1 = \text{Relu}(Z_1)$$

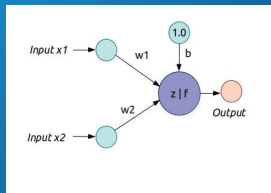


11/30/2024

pra-sâmi

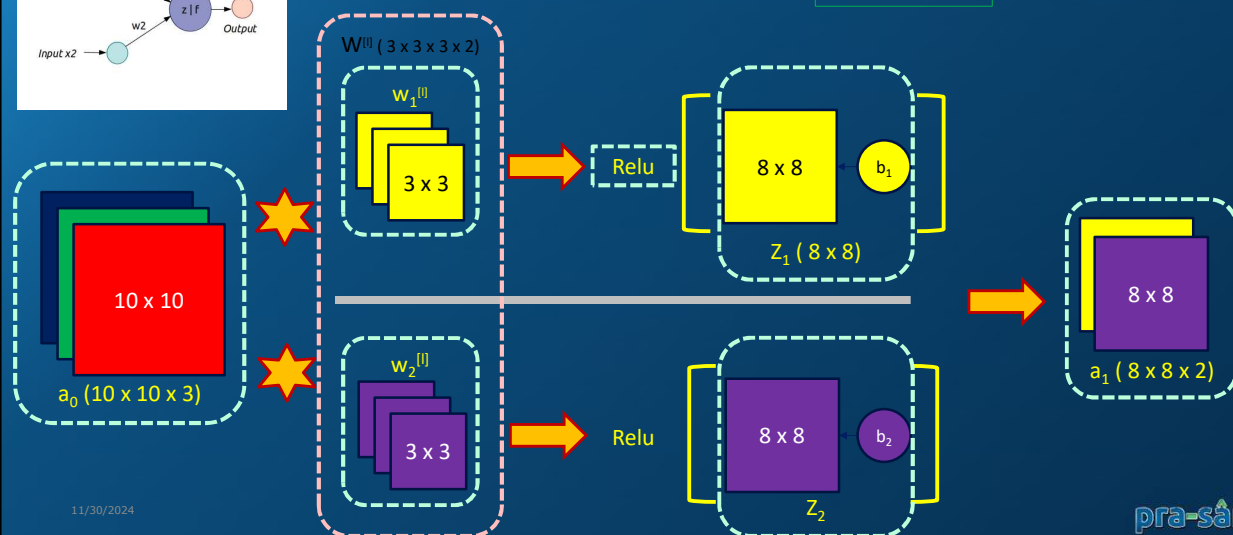
76

## One Layer of Conv Net



$$Z_1 = a_0 \cdot W_1 + b_1$$

$$a_1 = \text{Relu}(Z_1)$$



11/30/2024

pra-sâmi



79

## Lets Look at the Dimensions...

$f^{[l]}$ :	Filter Size	Input:	$n^{[l-1]}_H \times n^{[l-1]}_W \times n^{[l-1]}_C$
$p^{[l]}$ :	Padding size	Output:	$n^{[l]}_H \times n^{[l]}_W \times n^{[l]}_C$
$s^{[l]}$ :	Stride	$n^{[l]}_H$ :	$(n^{[l-1]}_H + 2 p^{[l]} - f^{[l]}) / s^{[l]} + 1$
$n^{[l]}_C$ :	Number of filters	$n^{[l]}_W$ :	$(n^{[l-1]}_W + 2 p^{[l]} - f^{[l]}) / s^{[l]} + 1$
Filter size:	$f^{[l]} \times f^{[l]} \times n^{[l-1]}_C$	Activations $a^{[l]}$ :	$n^{[l]}_H \times n^{[l]}_W \times n^{[l]}_C$
Weights (all filters):	$f^{[l]} \times f^{[l]} \times n^{[l-1]}_C \times n^{[l]}_C$	Biases:	$n^{[l]}_C$

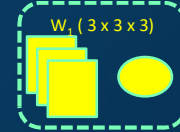
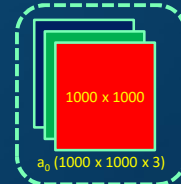
Weights are tensors of rank 4

Activation for all  $m$  training examples  $m$

$m \times n^{[l]}_H \times n^{[l]}_W \times n^{[l]}_C$

Don't be surprised if you see Filter number first

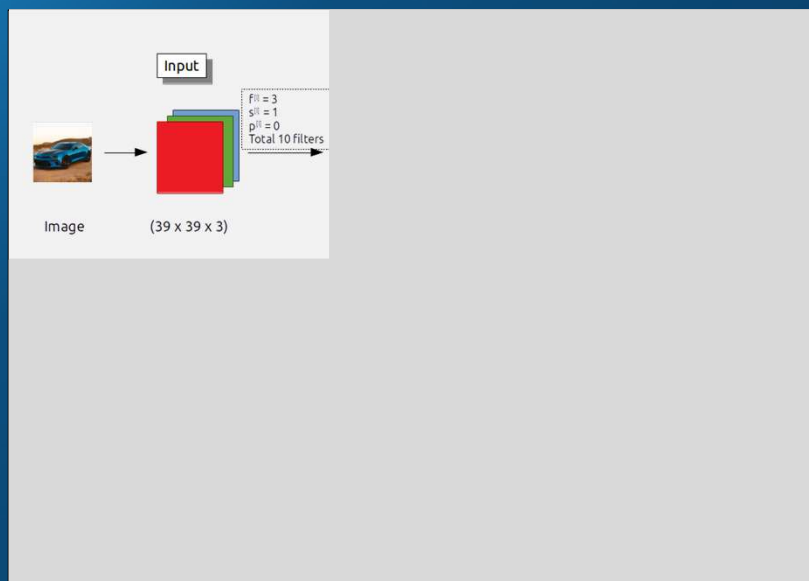
11/30/2024



pra-sami

80

## A Simple CNN with Conv Layers



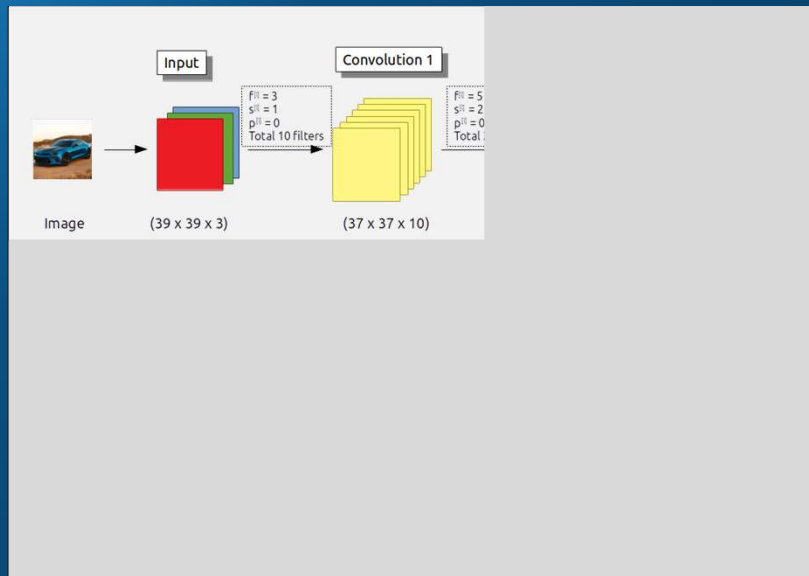
11/30/2024

pra-sami



81

## A Simple CNN with Conv Layers

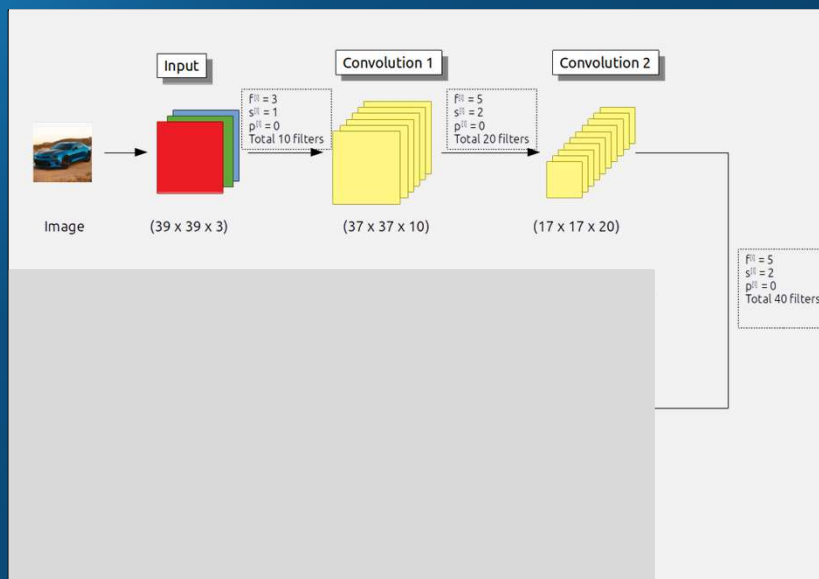


11/30/2024

pra-sami

82

## A Simple CNN with Conv Layers

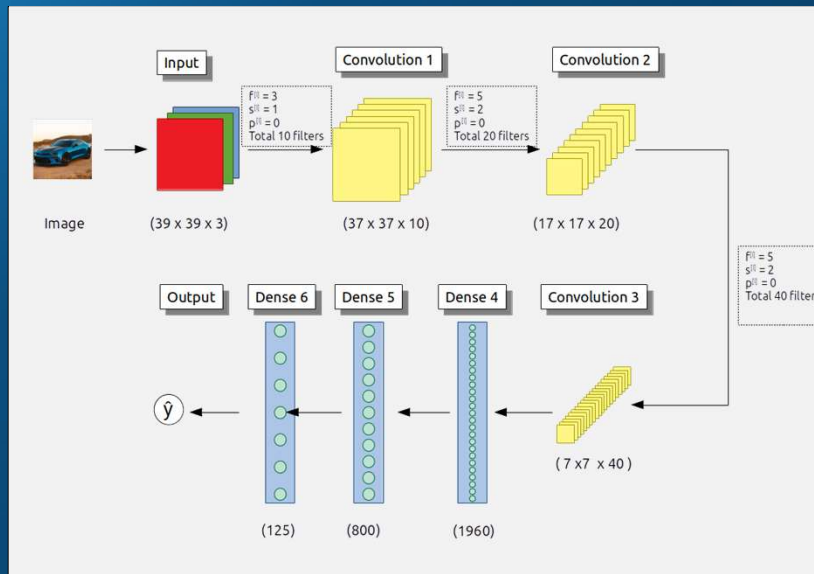


11/30/2024

pra-sami

83

## A Simple CNN with Conv Layers



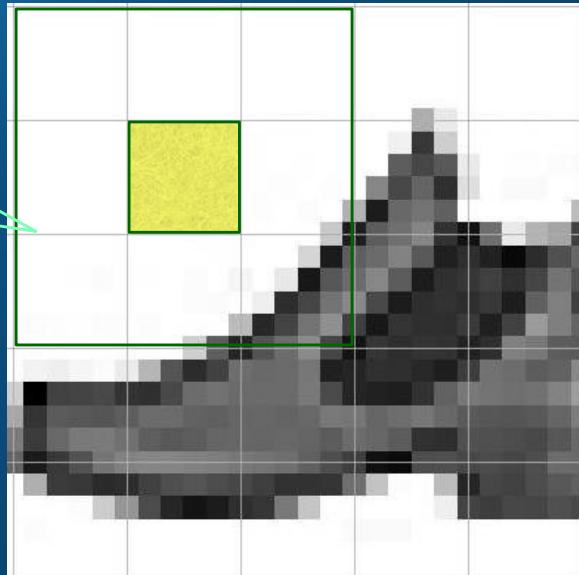
11/30/2024

pra-samí

84

## Convolution – Applying Filters

*9 datapoints  
result in one*



11/30/2024

pra-samí

85

Pooling...  
What is most significant in this area...

11/30/2024

pra-sami

86

## Pooling

- Two methods of Pooling – 'Max' and 'Average'
- Max : maximum value of the from the cells being filtered
- Average : Average Values from the cells

0	0	0	225	225	0	0	0
0	0	0	225	225	0	0	0
0	0	0	225	225	0	0	0
0	0	0	225	225	0	0	0
0	0	0	225	225	0	0	0
0	0	0	225	225	0	0	0
0	0	0	225	225	0	0	0
0	0	0	225	225	0	0	0

11/30/2024

- Mode = 'max'; pool = 2; stride = 2

0	225	225	0
0	225	225	0
0	225	225	0
0	225	225	0

pra-sami

87

## Other image

- Image Size = 10,10,3; filter Size = 3,3,3; stride = 1
- After Convolution = 8, 8, 1



8	-67	23	43	55	-72	-12	-30	81
7	313	343	0	0	0	0	-362	-291
6	559	390	0	0	0	0	-423	-601
5	498	390	0	0	0	0	-423	-633
4	498	390	0	0	0	0	-423	-633
3	559	390	0	0	0	0	-423	-601
2	318	344	0	0	0	0	-367	-296
1	-62	24	43	55	-72	-12	-35	76
0								

- Input size = 8,8,1; pool = 2; Stride = 2
- After pooling = 4,4,1

4	343	55	0	81
3	559	0	0	-423
2	559	0	0	-423
1	344	55	0	76
0				

11/30/2024

pra-sami

88

## Pooling

- Image Size = 10,10,3; filter Size = 3,3,3; stride = 1
- After Convolution = 8, 8, 1

8	-67	23	43	55	-72	-12	-30	81
7	313	343	0	0	0	0	-362	-291
6	559	390	0	0	0	0	-423	-601
5	498	390	0	0	0	0	-423	-633
4	498	390	0	0	0	0	-423	-633
3	559	390	0	0	0	0	-423	-601
2	318	344	0	0	0	0	-367	-296
1	-62	24	43	55	-72	-12	-35	76
0								

- Input size = 8,8,1; pool = 2; Stride = 2
- After pooling = 4,4,1
- Formula for size are still applicable,
- Its independently done on each channels
- Other option is to use Average instead of Max
  - ❖ But not used frequently.

4	343	55	0	81
3	559	0	0	-423
2	559	0	0	-423
1	344	55	0	76
0				

11/30/2024

pra-sami

89

## Pooling

- Image Size = 10,10,3; filter Size = 3,3,3; stride = 1
- After Convolution = 8, 8, 1

- Input size = 8,8,1; pool = 2; Stride = 2
- After pooling = 4,4,1
- Formula for size are still applicable,
- Its independently done on each channels
- Other is Average as expected but not used

Consider that each area represents presence of some feature in the image and high number represents, presence of that feature...

It has three (mode, pool and stride) hyperparameters to tune...

but no parameters to learn...

Gradient descent is not going to do anything here.... ☺

8									
7	313	343	0						
6									
5									
4	498	390	0	0	0	0	-423	-553	
3	498	390	0	0	0	0	-423	-553	
2	559	390	0	0	0	0	-423	-601	
1	318	344	0	0	0	0	-367	-296	
0	-62	24	43	55	-72	-12	-35	76	

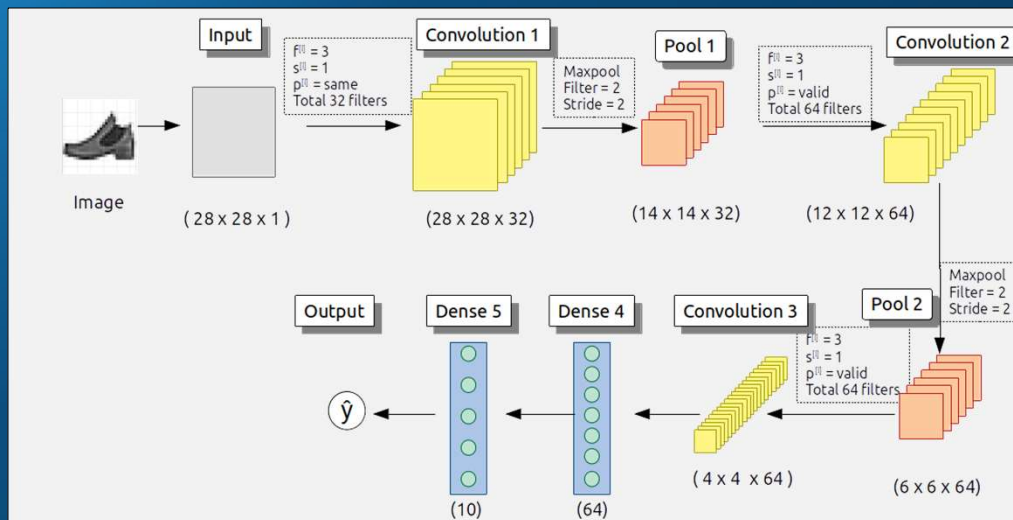
4									
3									
2									
1									
0									

11/30/2024

pra-sami

90

## Demo Example – Fashion MNIST

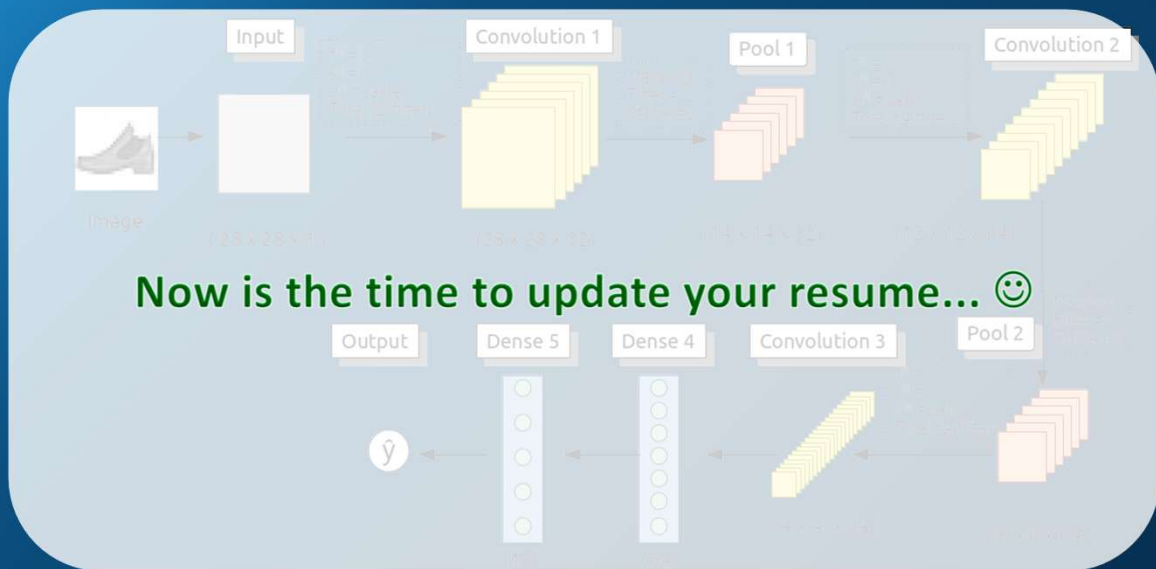


11/30/2024

pra-sami

91

Congratulations!!!!



11/30/2024

pra-sâmi

92

Reflect...

- ❑ What is the main purpose of Convolutional Neural Networks (CNNs)?
  - ❖ A) Text generation
  - ❖ B) Image and video recognition
  - ❖ C) Time series analysis
  - ❖ D) Natural language processing
- ❑ Answer: B) Image and video recognition
- ❑ Which of the following layers is a key component of CNNs?
  - ❖ A) Recurrent layer
  - ❖ B) Convolutional layer
  - ❖ C) Dropout layer
  - ❖ D) Activation layer
- ❑ Answer: B) Convolutional layer
- ❑ What is the role of the pooling layer in a CNN?
  - ❖ A) To increase the dimensions of the input
  - ❖ B) To extract features from the input data
  - ❖ C) To reduce the spatial dimensions (width and height) of the input
  - ❖ D) To combine multiple input channels
- ❑ Answer: C) To reduce the spatial dimensions (width and height) of the input
- ❑ Which operation is performed by a convolutional layer in a CNN?
  - ❖ A) Element-wise addition of the input matrix
  - ❖ B) Cross Correlation between a filter (kernel) and portions of the input matrix
  - ❖ C) Subtraction of one input channel from another
  - ❖ D) Matrix inversion
- ❑ Answer: B) Cross Correlation between a filter (kernel) and portions of the input matrix

11/30/2024

pra-sâmi



93

## Reflect...

- ❑ What is the function of a kernel (filter) in a CNN?
  - ❖ A) To resize images
  - ❖ B) To detect specific features like edges or textures in the input
  - ❖ C) To add noise to the image
  - ❖ D) To combine multiple images into one
- ❑ Answer: B) To detect specific features like edges or textures in the input
- ❑ What does stride refer to in a CNN?
  - ❖ A) The number of filters used
  - ❖ B) The number of steps the filter moves across the input matrix
  - ❖ C) The size of the input image
  - ❖ D) The number of output channels
- ❑ Answer: B) The number of steps the filter moves across the input matrix

- ❑ Which of the following is a common activation function used in CNNs?
  - ❖ A) Sigmoid
  - ❖ B) Tanh
  - ❖ C) ReLU (Rectified Linear Unit)
  - ❖ D) SoftMax
- ❑ Answer: C) ReLU (Rectified Linear Unit)
- ❑ In CNNs, what is the effect of padding?
  - ❖ A) To increase the number of filters
  - ❖ B) To prevent the reduction of spatial dimensions by adding zeros around the input matrix
  - ❖ C) To reduce the memory footprint of the model
  - ❖ D) To change the size of the kernel
- ❑ Answer: B) To prevent the reduction of spatial dimensions by adding zeros around the input matrix

11/30/2024

pra-sâmi

94

## Reflect...

- ❑ Why do deeper CNNs typically perform better than shallow CNNs?
  - ❖ A) Deeper CNNs have more parameters and can memorize the training data better
  - ❖ B) Deeper CNNs can learn hierarchical representations, capturing complex features
  - ❖ C) Deeper CNNs require fewer data points for training
  - ❖ D) Deeper CNNs prevent overfitting
- ❑ Answer: B) Deeper CNNs can learn hierarchical representations, capturing complex features
- ❑ What is the vanishing gradient problem in the context of CNNs?
  - ❖ A) The gradients become too small to update the weights effectively in deeper networks
  - ❖ B) The model loses information about small details in images
  - ❖ C) The network stops learning after a certain number of epochs
  - ❖ D) The gradients become too large, leading to unstable training
- ❑ Answer: A) The gradients become too small to update the weights effectively in deeper networks

- ❑ What is the purpose of the fully connected (dense) layer at the end of a CNN?
  - ❖ A) To downsample the input
  - ❖ B) To map the learned features to the final output classes
  - ❖ C) To combine the pooling layers into a single layer
  - ❖ D) To prevent overfitting by regularizing the model
- ❑ Answer: B) To map the learned features to the final output classes
- ❑ What is a common method to reduce overfitting in a CNN?
  - ❖ A) Use a very large filter size
  - ❖ B) Add more fully connected layers
  - ❖ C) Use techniques like dropout, data augmentation, or early stopping
  - ❖ D) Increase the learning rate
- ❑ Answer: C) Use techniques like dropout, data augmentation, or early stopping

11/30/2024

pra-sâmi

95

**THANK YOU**

11/30/2024

pra-samí