



CONVOLUTION NEURAL NETWORKS – CNN PART I

Deep Neural Networks

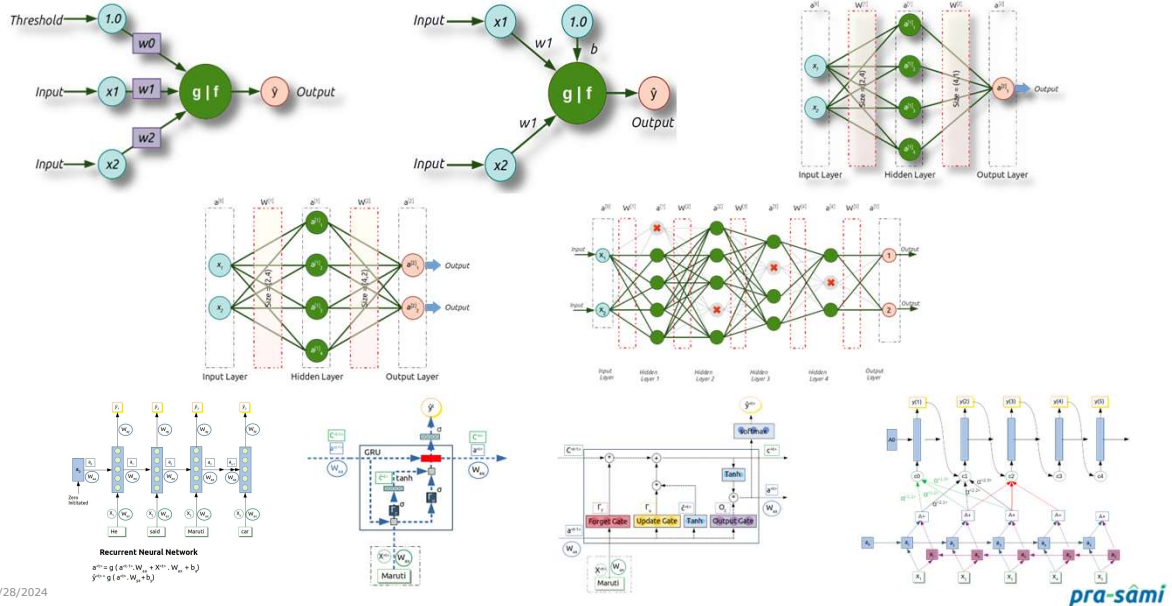
Session 21

Pramod Sharma

pramod.sharma@prasami.com

2

Story So Far...



3

Agenda



5/28/2024

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.

5/28/2024

pra-sâmi

5

What is Computer Vision...

5/28/2024

pra-sâmi

6

Ambulance given green light all through....

5/28/2024

pra-sâmi

7

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

5/28/2024

pra-sâmi

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

5/28/2024

pra-sâmi

9

Computer Vision

- ❑ Self driving car
- ❑ Fully automated warehouse and ports
 - ❖ <https://youtu.be/RFV8IkY52iY>
- ❑ 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

5/28/2024

pra-sâmi

10

Computer Vision - Style transfer

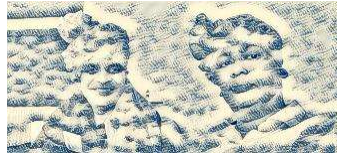
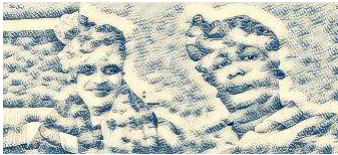
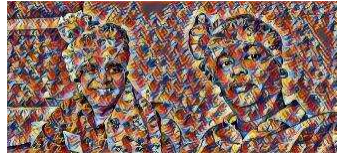


5/28/2024

pra-sâmi

11

Computer Vision - Style transfer



5/28/2024

pra-sâmi

12

Computer Vision - Style transfer



5/28/2024

pra-sâmi

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)

5/28/2024

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

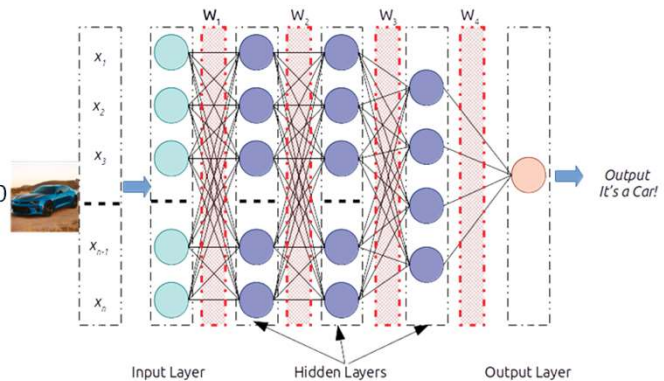
5/28/2024

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



5/28/2024

pra-sâmi

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?

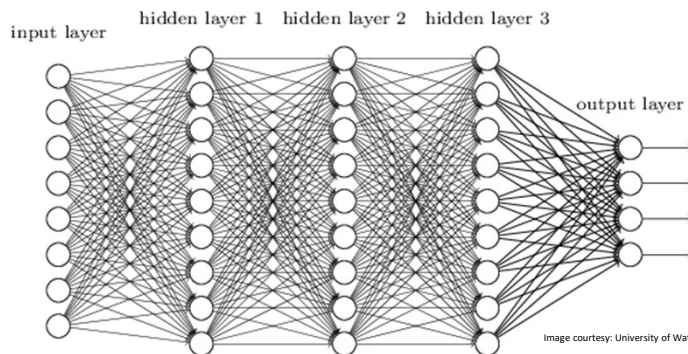


Image courtesy: University of Waterloo.

5/28/2024

pra-sâmi

17

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

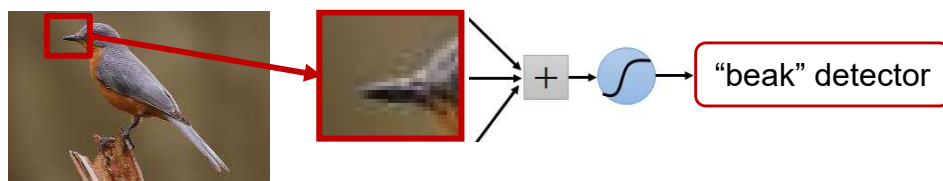
5/28/2024

pra-sâmi

18

Learning an image...

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



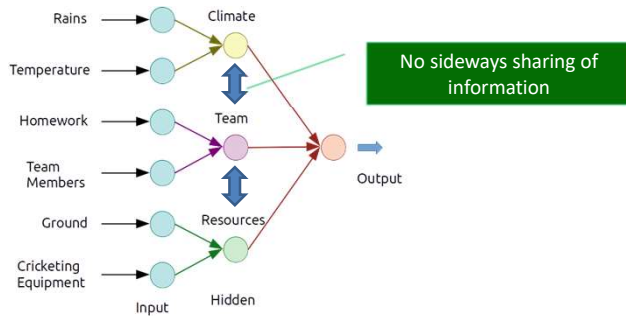
5/28/2024

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”

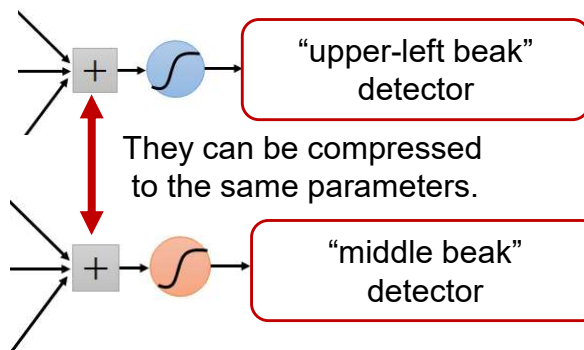


5/28/2024

pra-sâmi

20

Learning an image...



5/28/2024

pra-sâmi

21

Learning an image...

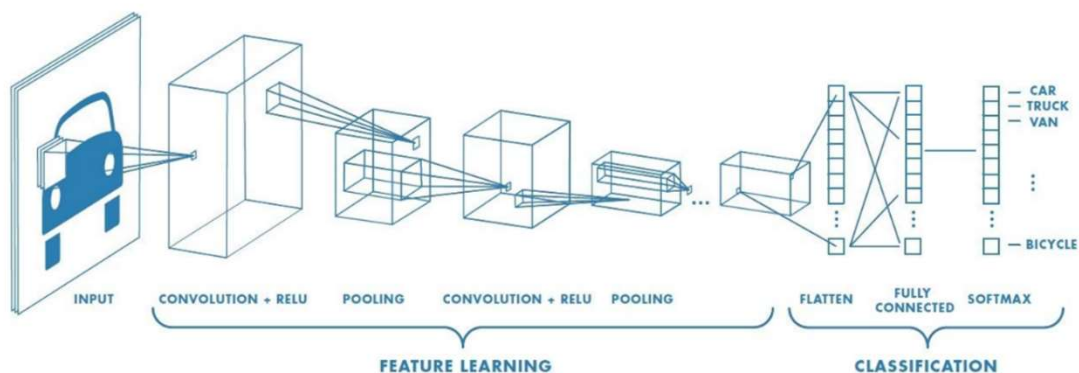
- ❑ The same sorts of features that are useful in analyzing one part of the image will probably be useful for analyzing other parts as well.
 - ❖ E.g., edges, corners, contours, object parts
- ❑ We want a neural net architecture that lets us learn a set of feature detectors that are applied at all image 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...

5/28/2024

pra-sâmi

22

Overall Layout

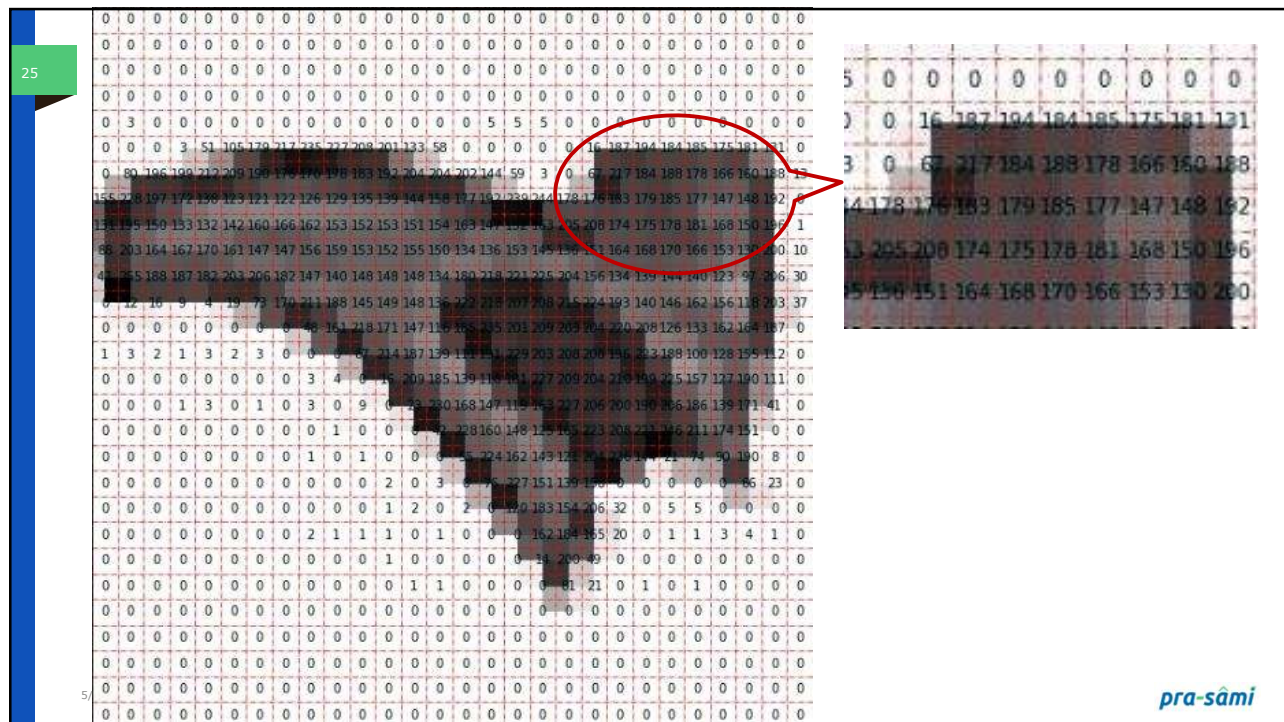


5/28/2024

pra-sâmi

pra-sâmi

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

Diagram illustrating a convolution operation. An input grid is processed by a filter (kernel) to produce an output grid. The filter is labeled "A filter" and the output is labeled "Outputs". A red arrow points to the filter, and a green arrow points to the output.

pra-sâmi

27

What does this Convolution Filter/ Kernel do?


 $*$

0	1	0
1	4	1
0	1	0



5/28/2024

pra-sâmi

28

What does this Convolution Filter/ Kernel do?


 $*$

0	-1	0
-1	8	-1
0	-1	0



5/28/2024

pra-sâmi

29

What does this Convolution Filter/ Kernel do?


 $*$

0	-1	0
-1	4	-1
0	-1	0



5/28/2024

pra-sâmi

30

What does this Convolution Filter/ Kernel do?


 $*$

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



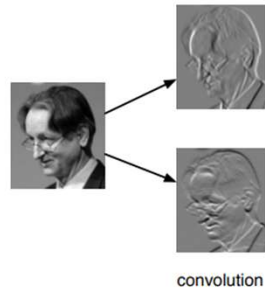
5/28/2024

pra-sâmi

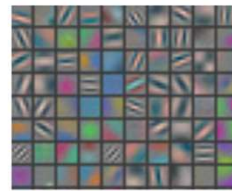
31

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.



Example first-layer filters



(Zeiler and Fergus, 2013, Visualizing and understanding convolutional networks)

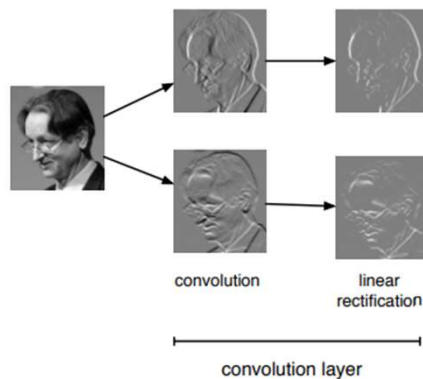
5/28/2024

pra-sâmi

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.



- ❑ 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

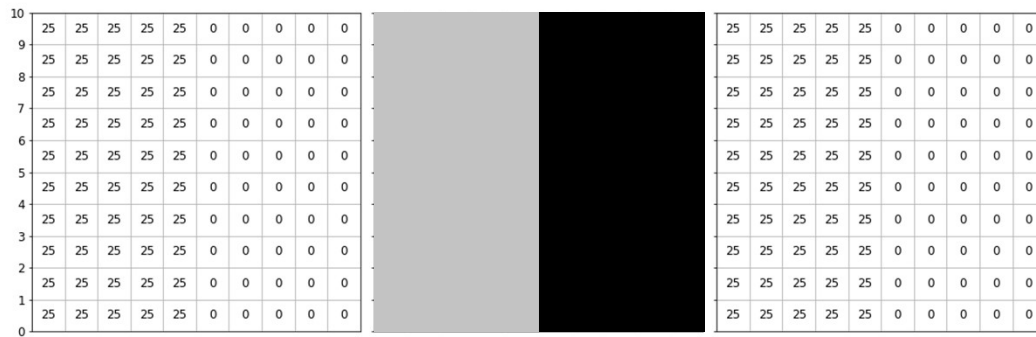
5/28

pra-sâmi

33

Image with a edge

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



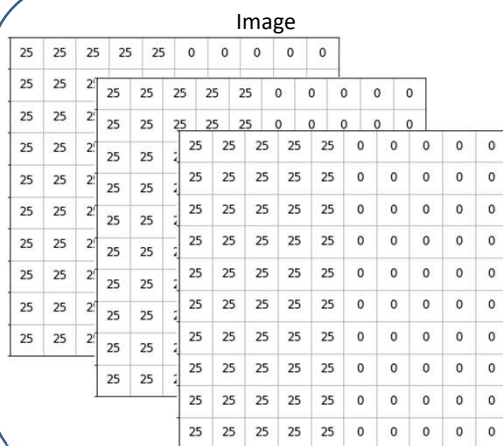
- Apply filters on the image!

5/28/2024

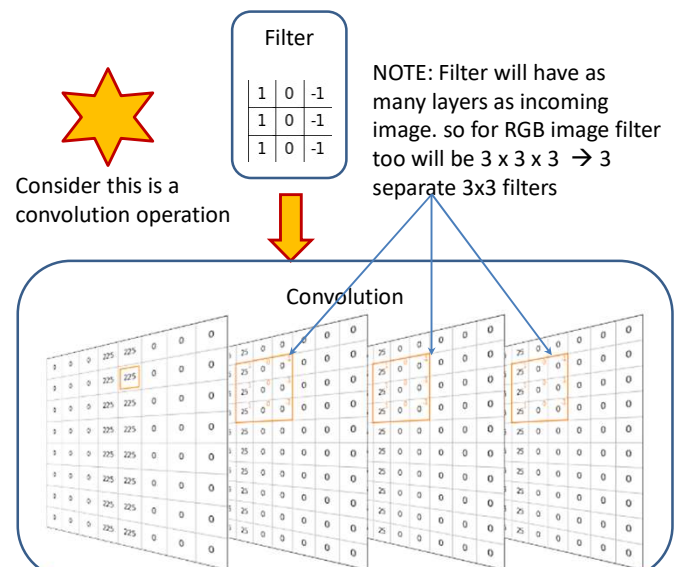
pra-sâmi

34

Convolution on 3D images (RGB)



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



5/28/2024

pra-sâmi

35

Convolution on 3D images (RGB)

□ First convolution

25	¹ 25	⁰ 25	⁻¹ 25	25	25	0	0	0	0	0
25	¹ 25	⁰ 25	⁻¹ 25	25	25	0	0	0	0	0
25	¹ 25	⁰ 25	⁻¹ 25	25	25	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0	0

5/28/2024

□ 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

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

$$\diamond = 0$$

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

pra-sâmi

36

Convolution on 3D images (RGB)

□ Second convolution

❖ It will be identical to First

25	¹ 25	⁰ 25	⁻¹ 25	25	0	0	0	0	0	0
25	¹ 25	⁰ 25	⁻¹ 25	25	0	0	0	0	0	0
25	¹ 25	⁰ 25	⁻¹ 25	25	0	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0	0
25	25	25	25	25	0	0	0	0	0	0

5/28/2024

□ 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

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

$$\diamond = 0$$

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

pra-sâmi

37

Convolution on 3D images (RGB)

□ What happens 4th step

25	25	25	25 ¹	25 ⁰	0 ⁻¹	0	0	0	0
25	25	25	25 ¹	25 ⁰	0 ⁻¹	0	0	0	0
25	25	25	25 ¹	25 ⁰	0 ⁻¹	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

5/28/2024

□ 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

pra-sâmi

38

Convolution on 3D images (RGB)

□ And for 5th Step

25	25	25	25	25 ¹	0 ⁰	0 ⁻¹	0	0	0
25	25	25	25	25 ¹	0 ⁰	0 ⁻¹	0	0	0
25	25	25	25	25 ¹	0 ⁰	0 ⁻¹	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

5/28/2024

□ 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

pra-sâmi

39

Convolution on 3D images (RGB)

□ 6th Step onwards again all values are 0

25	25	25	25	25	0 ¹	0 ⁰	0 ⁻¹	0	0
25	25	25	25	25	0 ¹	0 ⁰	0 ⁻¹	0	0
25	25	25	25	25	0 ¹	0 ⁰	0 ⁻¹	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

5/28/2024

pra-sâmi

40

Convolution

Image

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



Filter

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



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

Convolution

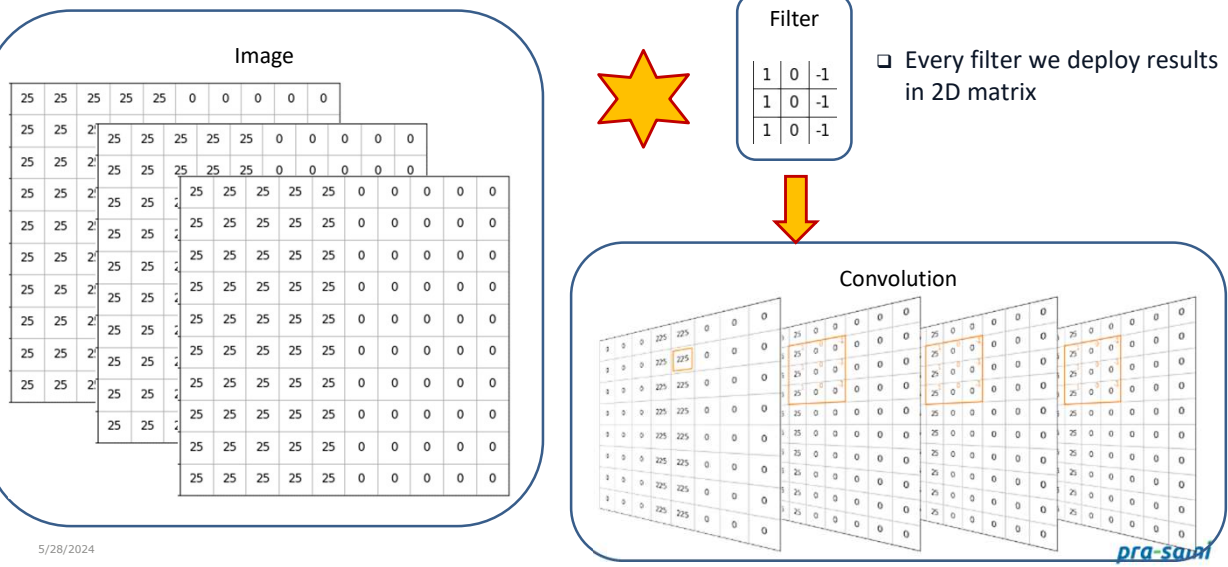
25	25	25	25	25	0	0	0	0	0
25	25	25	25	25	25	0	0	0	0
25	25	25	25	25	25	0	0	0	0
25	25	25	25	25	25	0	0	0	0
25	25	25	25	25	25	0	0	0	0
25	25	25	25	25	25	0	0	0	0
25	25	25	25	25	25	0	0	0	0
25	25	25	25	25	25	0	0	0	0
25	25	25	25	25	25	0	0	0	0
25	25	25	25	25	25	0	0	0	0

5/28/2024

pra-sâmi

41

Convolution



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

5/28/2024

pra-sâmi

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

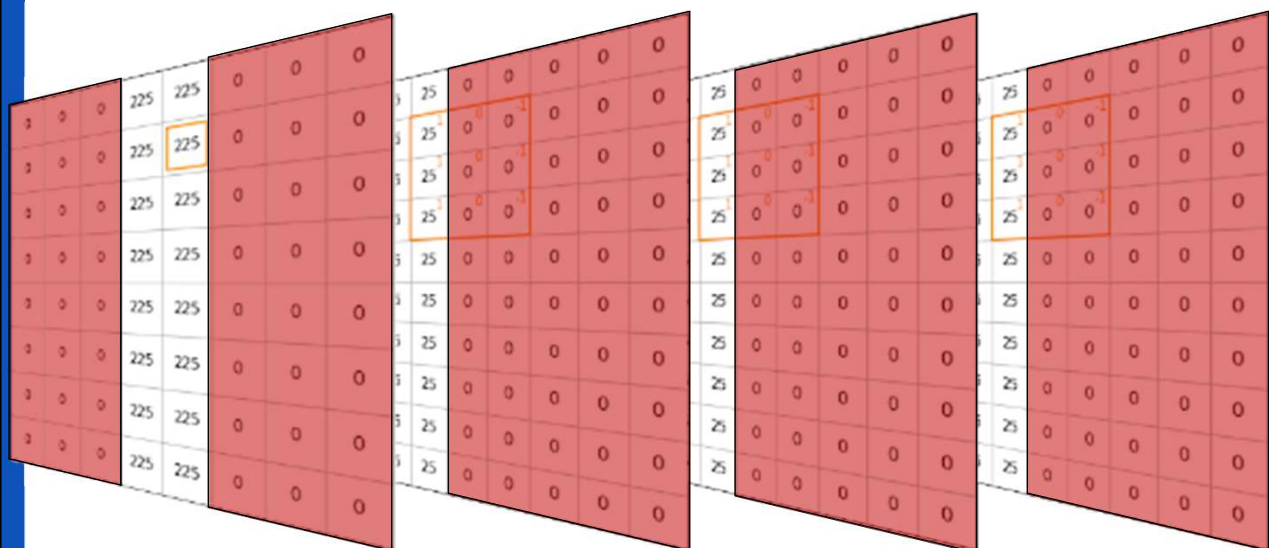
□ Hence, Output image size
 $= \{nH_{in} - nF + 1\} \times \{nW_{in} - nF + 1\}$

5/28/2024

pra-sâmi

44

Has it detected the edge?



5/28/2024

pra-sâmi

45

What if we move two steps at a time

5/28/2024

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

5/28/2024

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

□ Hence, Output image size =

$$\{ (nH_{in} - nF) / stride + 1 \} \times \{ (nW_{in} - nF) / stride + 1 \}$$

□ If we apply multiple filters → this layer will have 3D matrix.

❖ Each layer corresponding to one filter.

5/28/2024

pra-sâmi

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

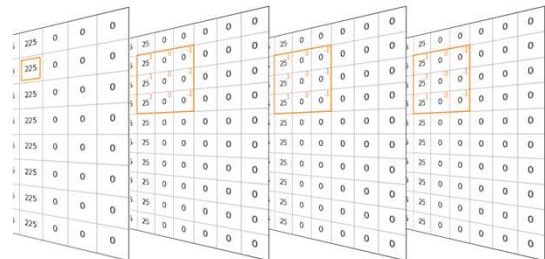
❖ one in most cases...

□ 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



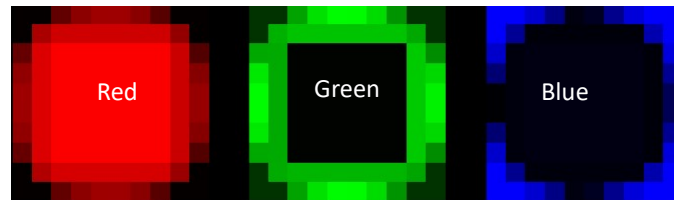
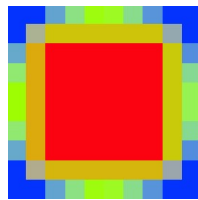
5/28/2024

pra-sâmi

49

Another image

- This image has a few distinct edges



10	3	3	98	140	158	157	137	87	3	3	51	51	163	225	252	250	220	148	51	51	252	252	207	131	15	34	140	218	252	252
9	3	170	206	206	206	206	206	206	148	3	51	178	196	196	196	196	196	196	169	51	252	147	10	10	10	10	10	175	252	
8	102	220	252	252	252	252	252	252	201	82	169	170	3	3	3	3	3	3	203	141	203	13	18	18	18	18	18	18	10	222
7	145	220	252	252	252	252	252	252	201	134	233	170	3	3	3	3	3	3	203	217	113	13	18	18	18	18	18	18	10	147
6	158	220	252	252	252	252	252	252	201	151	252	170	3	3	3	3	3	3	203	242	3	13	18	18	18	18	18	18	10	84
5	158	220	252	252	252	252	252	252	201	151	252	170	3	3	3	3	3	3	203	242	3	13	18	18	18	18	18	18	10	84
4	145	220	252	252	252	252	252	252	201	134	233	170	3	3	3	3	3	3	203	217	113	13	18	18	18	18	18	18	10	147
3	102	220	252	252	252	252	252	252	201	82	169	170	3	3	3	3	3	3	203	141	203	13	18	18	18	18	18	18	10	222
2	3	177	213	213	213	213	213	213	154	3	51	172	182	182	182	182	182	182	166	51	252	142	12	12	12	12	12	12	172	252
1	3	3	98	140	158	157	137	87	3	3	51	51	163	225	252	250	220	148	51	51	252	252	207	131	15	34	140	218	252	252
0																														

5/28/2024

pra-sâmi

50

Another Convolution

Image

3	3	98	51	51	163	225	252	250	220	148	51	51			
3	170	206	51	178	196	252	252	207	131	15	34	140	218	252	252
102	220	252	169	170	3	252	147	10	10	10	10	10	10	175	252
145	220	252	233	170	3	203	13	18	18	18	18	18	18	10	222
158	220	252	252	170	3	113	13	18	18	18	18	18	18	10	147
158	220	252	252	170	3	3	13	18	18	18	18	18	18	10	84
145	220	252	233	170	3	3	13	18	18	18	18	18	18	10	84
102	220	252	233	170	3	113	13	18	18	18	18	18	18	10	147
3	177	213	169	170	3	203	13	18	18	18	18	18	18	10	222
3	3	98	51	172	182	252	142	12	12	12	12	12	12	172	252
			51	51	163	252	252	207	131	15	34	140	218	252	252



Filter

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



Convolution

3	3	98	51	51	163	225	252	252	207	131	-67	23	43	55	-72	-12	-18	81
3	170	206	51	178	196	196	252	147	10	10	313	343	0	0	0	0	362	298
102	220	252	169	170	3	3	203	13	18	18	559	390	0	0	0	0	433	433
145	220	252	233	170	3	3	113	13	18	18	498	390	0	0	0	0	433	433
158	220	252	252	170	3	3	3	13	18	18	498	390	0	0	0	0	433	433
158	220	252	252	170	3	3	3	13	18	18	559	390	0	0	0	0	433	433
145	220	252	233	170	3	3	113	13	18	18	318	344	0	0	0	0	367	298
102	220	252	169	170	3	3	203	13	18	18	-62	24	43	55	-72	-12	-18	81
3	177	213	51	172	182	225	252	142	12	12								
3	3	98	51	51	163	225	252	252	207	131								

5/28/2024

pra-sâmi

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 = $-448 + -91 + 472 = -67$

5/28/2024

pra-sâmi

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

5/28/2024

pra-sâmi

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 = $-448 + -91 + 472 = -67$

5/28/2024

- 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

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

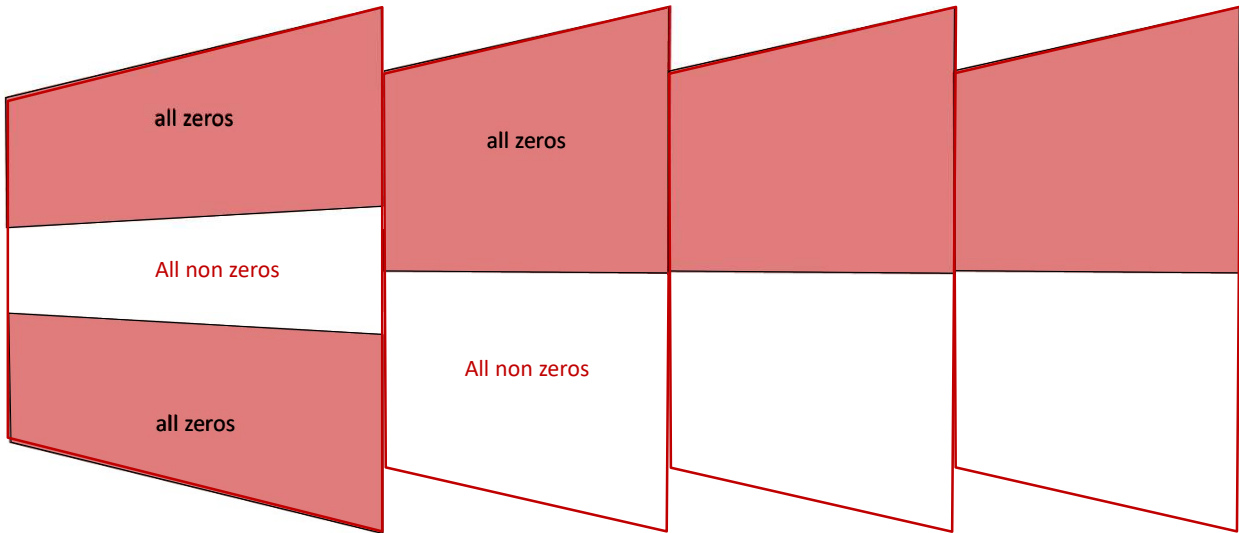
- The math will be exactly the same and we would get horizontal edge

5/28/2024

pra-sâmi

55

Horizontal Edge...



5/28/2024

pra-sâmi

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

5/28/2024

pra-sâmi

57

Shade Reversal

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

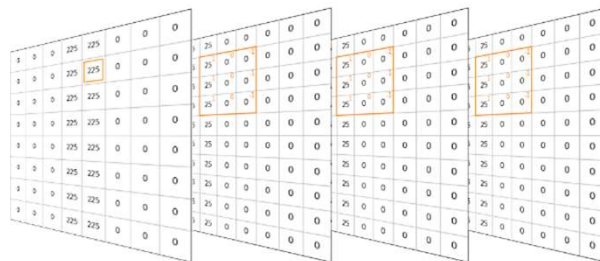
5/28/2024

pra-sâmi

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...



5/28/2024

pra-sâmi

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...

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

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

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

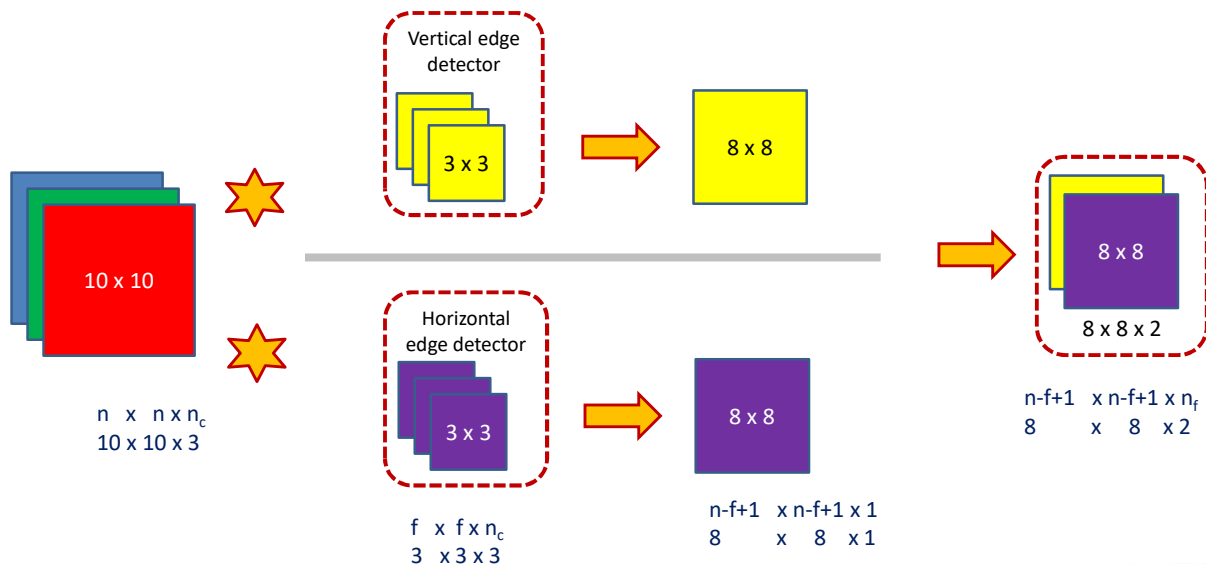
- So far we have been showing that 3D image converts to 2D image when we apply filter
- By applying a number of filters to detect different edges, we can have 3d Convolutional Volumes.

5/28/2024

pra-sâmi

60

Multiple Filters



5/28/2024

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

5/28/2024

pra-sâmi

62

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

5/28/2024

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

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

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

We are back to original size...

5/28/2024

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

5/28/2024

pra-sâmi

67

How much to pad???

- ❑ With $p = (nF - 1)/2$ for 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

5/28/2024

pra-sâmi

68

Cross-Correlation vs. Convolution

5/28/2024

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

5/28/2024

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... 😊

5/28/2024

pra-sâmi

71

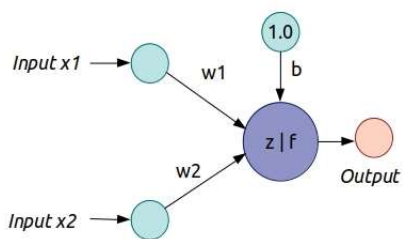
One layer of Convolutional Net

5/28/2024

pra-sâmi

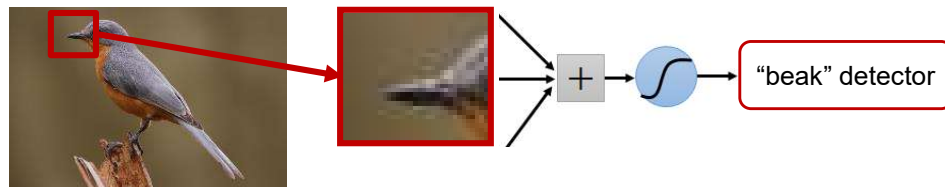
72

One layer of Conv Net



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

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

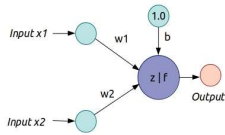


5/28/2024

pra-sâmi

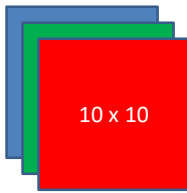
73

One Layer of Conv Net



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

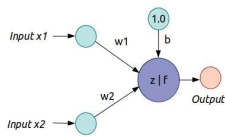
$$a_1 = \text{Relu}(Z_1)$$



5/28/2024

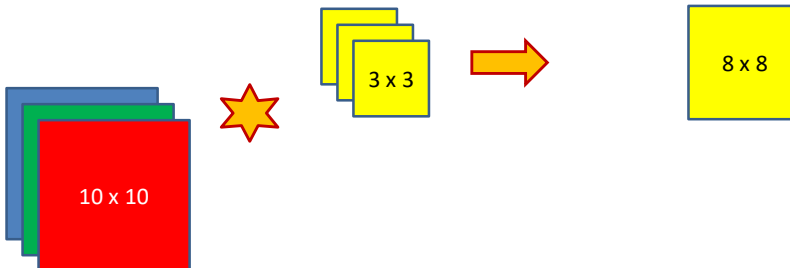
74

One Layer of Conv Net



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

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

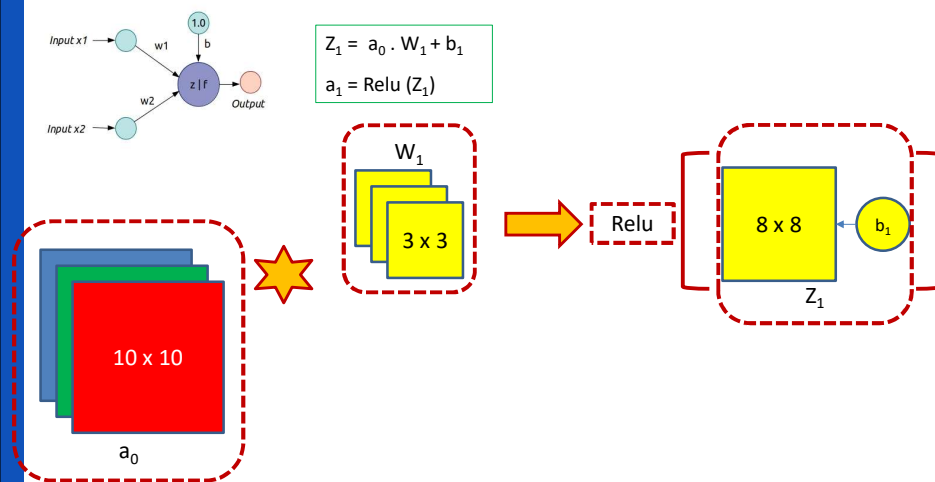


5/28/2024

pra-sami

75

One Layer of Conv Net

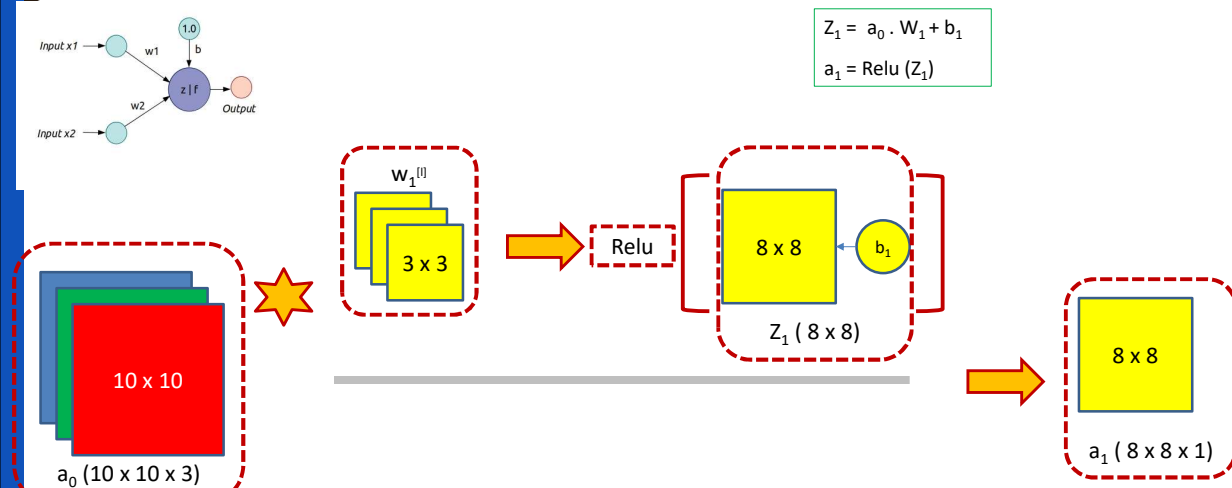


5/28/2024

pra-sami

76

One Layer of Conv Net

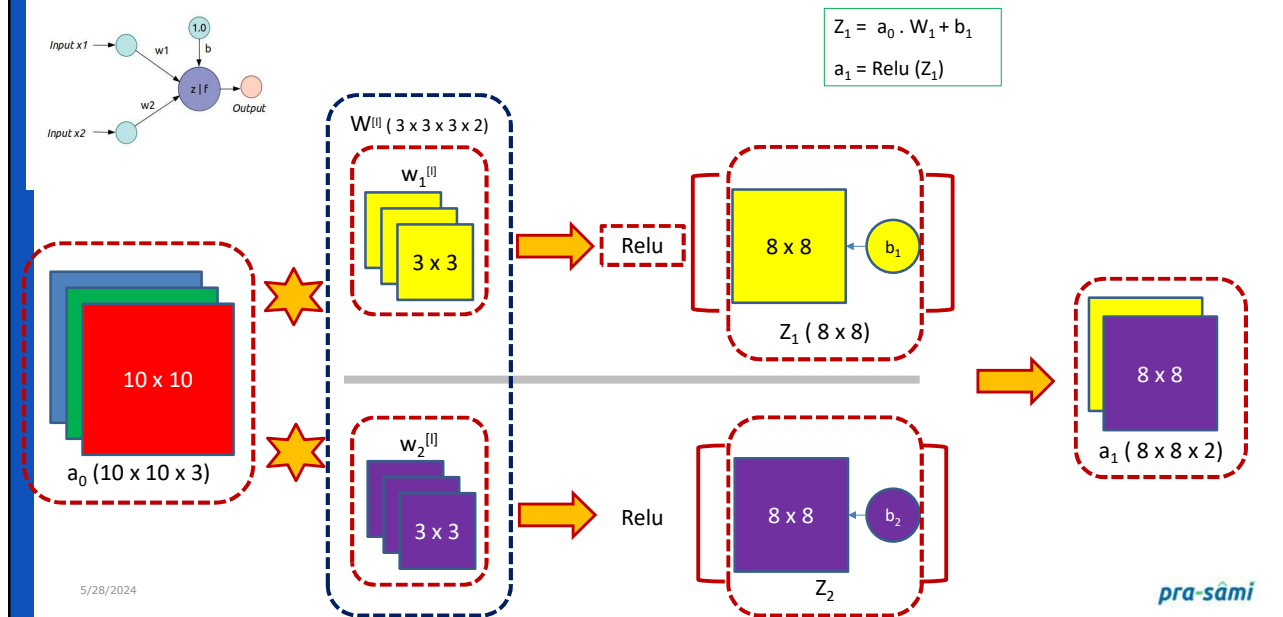


5/28/2024

pra-sami

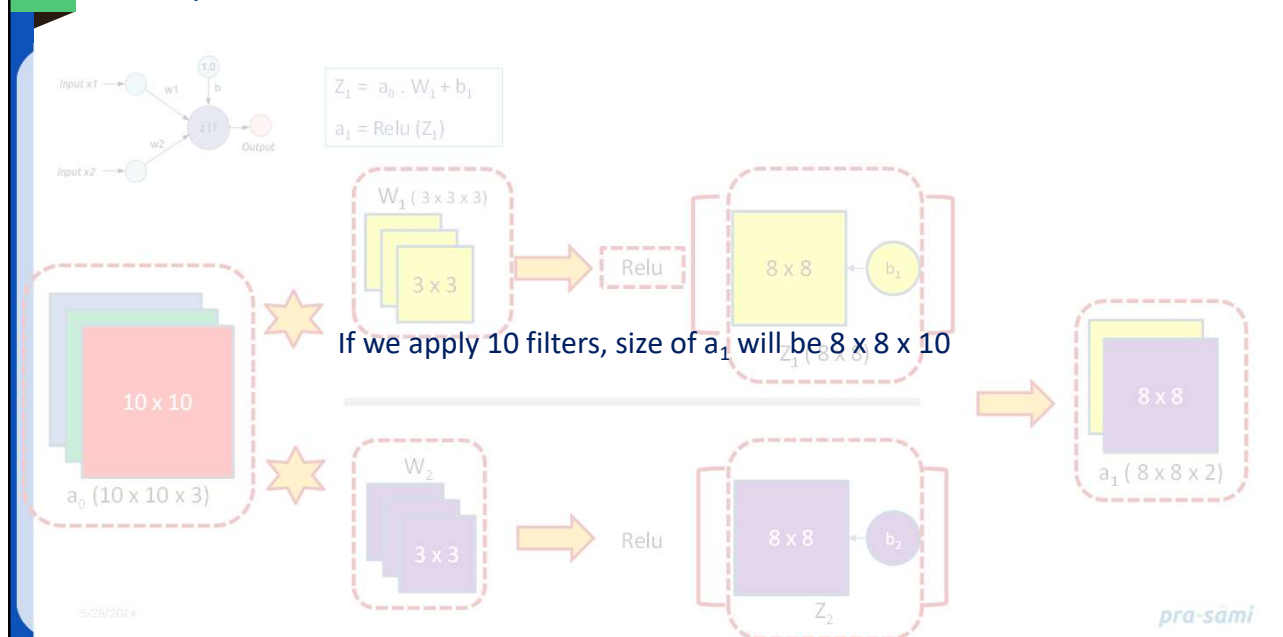
77

One Layer of Conv Net



78

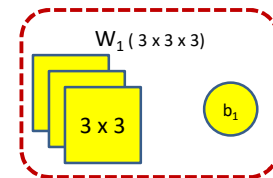
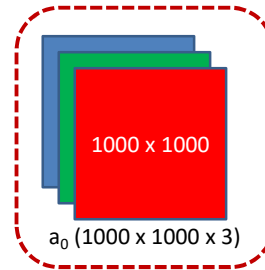
One Layer of Conv Net



79

How Many Parameters...

- ❑ Imagine using 10 filters in a layer, how many parameters in the layer
- ❑ Each filter is $3 \times 3 \times 3$ + a bias = 28 parameters
- ❑ For 10 filters = total 280 parameters
- ❑ Hence, irrespective of size of input image /activation, we still have 280 filters learning what we need to know.... Yay!!!!
- ❑ Helps in prevention of over-fitting



5/28/2024

pra-sâmi

80

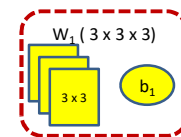
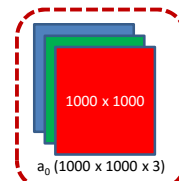
Lets Look at the Dimensions...

$f^{[l]}$:	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

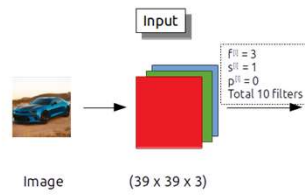


5/28/2024

pra-sâmi

81

A Simple CNN with Conv Layers

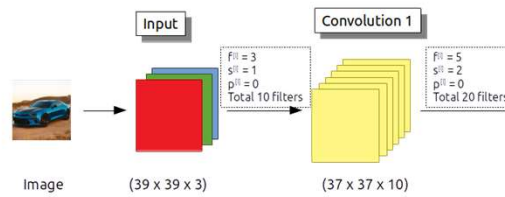


5/28/2024

pra-sâmi

82

A Simple CNN with Conv Layers

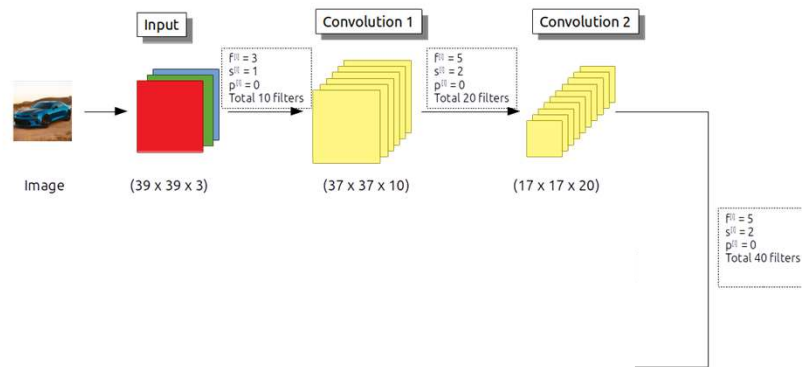


5/28/2024

pra-sâmi

83

A Simple CNN with Conv Layers

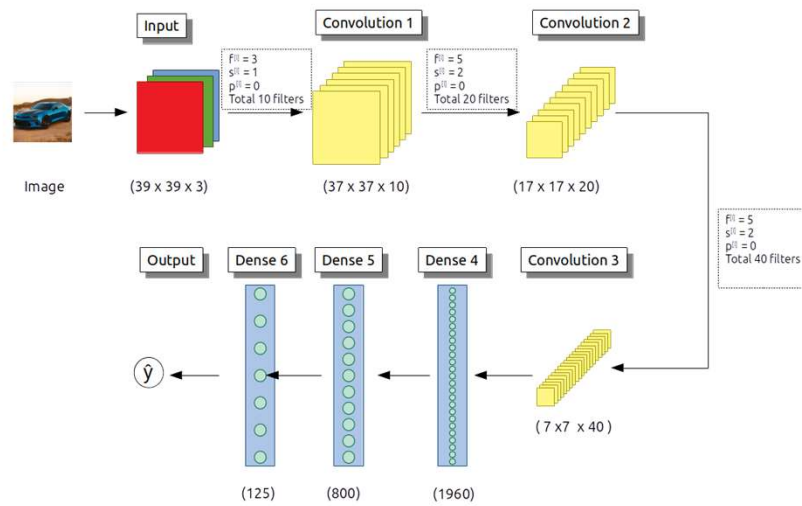


5/28/2024

pra-sâmi

84

A Simple CNN with Conv Layers



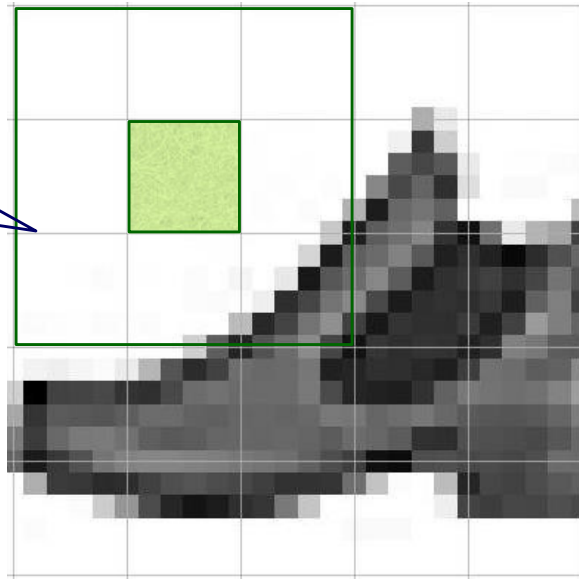
5/28/2024

pra-sâmi

85

Convolution – Applying Filters

*9 datapoints
result in one*



5/28/2024

pra-sâmi

86

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

5/28/2024

pra-sâmi

87

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

5/28/2024

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

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

pra-sâmi

88

Other image

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

-67	23	43	55	-72	-12	-30	81
313	343	0	0	0	0	-362	-291
559	390	0	0	0	0	-423	-601
498	390	0	0	0	0	-423	-633
498	390	0	0	0	0	-423	-633
559	390	0	0	0	0	-423	-601
318	344	0	0	0	0	-367	-296
-62	24	43	55	-72	-12	-35	76

5/28/2024

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

343	55	0	81
559	0	0	-423
559	0	0	-423
344	55	0	76

pra-sâmi

89

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				

5/28/2024

pra-sâmi

90

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	-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								

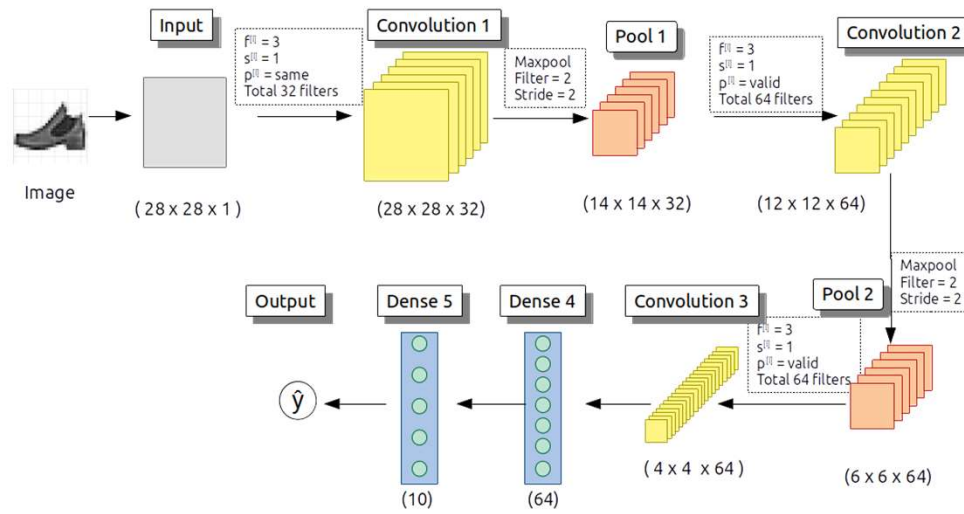
4	343	55	0	81
3	559	0	0	-423
2	559	0	0	-423
1	344	55	0	76
0				

5/28/2024

pra-sâmi

91

Demo Example – Fashion MNIST

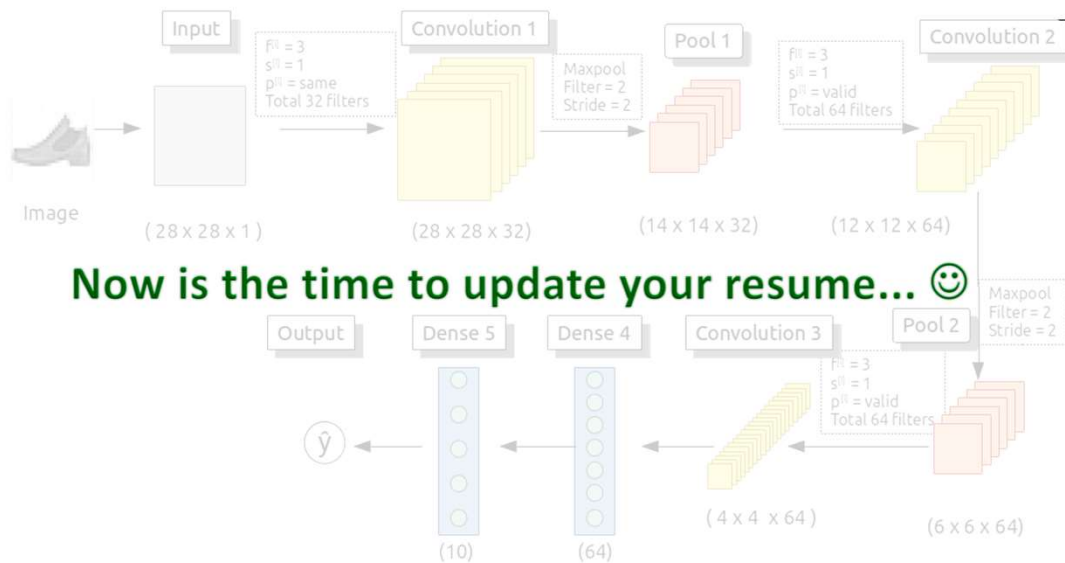


5/28/2024

pra-sâmi

92

Congratulations!!!!



5/28/2024

pra-sâmi

