

Convolution layers

Only a subset of I/P is connected to O/P.
Each O/P captures only a SMALL SET of adjacent input.

V.
local view
local patterns captured.
global view
Every O/P is connected to every I/P pattern.
O/P captures full I/P pattern.

Layers in MLP.

Many convolution layers capture different local patterns of I/P based on kernel.

Since "same kernel" slides over entire I/P,
Weight to be learnt <<< MLP.

1D Convolution

Input: 1D array

Kernel: Small 1D array

Output: 1D array

Types of filters.

No. of trainable params is HUGE.

Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	

also: filter — small matrix to apply effects to image

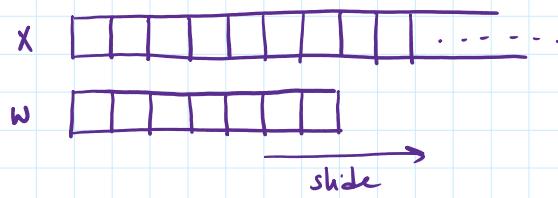
- ① Kernel slides over I/P data
- ② Elementwise multiplication (kernel \times I/P values)
- ③ Summation of multiplications — get single value
- ④ Repeat across entire I/P.

Stride: number of units the filter shifts @ each step

Padding → Valid: no padding — filter slides within bounds of I/P
smaller output — filter can't go beyond edges

Zero: extra zeroes added around edge of input
output size = input size.

Assume kernel has dimension k_w .



Output when you slide:

$$Y_n = \sum_{j=0}^{k_w} X_{n+j} \cdot w_j$$

Express this as matrix multiplication:

$$\begin{bmatrix} w_0 & w_1 & w_2 \end{bmatrix} \otimes \begin{bmatrix} x_0 & x_1 & x_2 & \dots & x_{n-1} & x_n \end{bmatrix}$$

$$\bar{w} \otimes \bar{x}$$

for stride=1, convolution weight matrix:

Extras:
for a CWM,
no. of rows = $\frac{(n - \text{size of kernel}) + 1}{\text{stride}}$
no. of columns = length of I/P vector (ie n)

$$\begin{bmatrix} w_0 & w_1 & w_2 & 0 & \dots & 0 & 0 & 0 \\ 0 & w_0 & w_1 & w_2 & \dots & 0 & 0 & 0 \\ \vdots & & & & & & & \\ 0 & 0 & 0 & 0 & \dots & w_0 & w_1 & w_2 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{bmatrix}$$

Convolution weight matrix Input matrix

→ if stride=2:

$$\begin{bmatrix} w_0 & w_1 & w_2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & w_0 & w_1 & w_2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & w_0 & w_1 & w_2 & 0 \end{bmatrix}$$

Independent of stride, kernel size.

16 n)

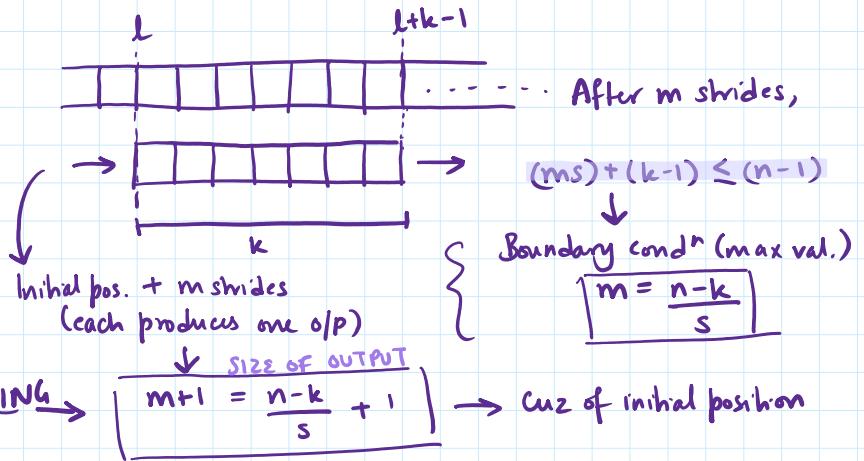
Independent of stride, kernel size.

Output size in 1D convolution:

$$|\text{Input}| = n$$

$$|\text{Kernel}| = k$$

$$|\text{Stride}| = s$$



If ZERO PADDING

0	0	0	0	0	0
0	35	19	25	6	0
0	13	22	16	53	0
0	4	3	7	10	0
0	9	8	1	3	0
0	0	0	0	0	0

$$\text{Input} = n + 2p \quad (\text{cuz } p \text{ zeroes on each side})$$

SIZE OF OUTPUT

$\text{Output} = n + 2p - k + 1$

*

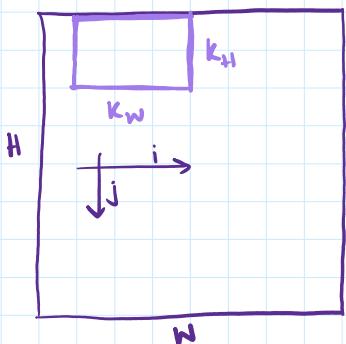
Blue Channel
Previous Patch (3x3)

Current Patch (3x3)



2 Dimensions Now.

→ Well just be sliding our tile on a wall.

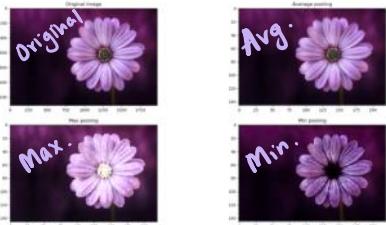


$$\rightarrow \text{now, } Y_{y,x} = \sum_j \sum_{i=0}^{K_H} X_{x+i, y+j} \cdot W_{ij}$$

Corner cases? $n \rightarrow n + S_n$
 $y \rightarrow y + S_H$

RHS top corner: $x = n$, $y = y + 1$
 RHS bottom corner: STOP.
 all ghost 1/P replaced by 0.

Types of Pooling



↑ If object in image moves:

Max pooling: capture most prominent features (OR)
Average -||- : smooth out variations.

POOLING → Max (max. value of each patch from feature map)
→ Average (same, but take avg.)

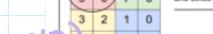
Downsampling → return same info →

Pooling \approx Local translation invariance.

→ Minor changes in I/P features \Rightarrow diff. O/P feature map
Now what?

less computation → less overfitting

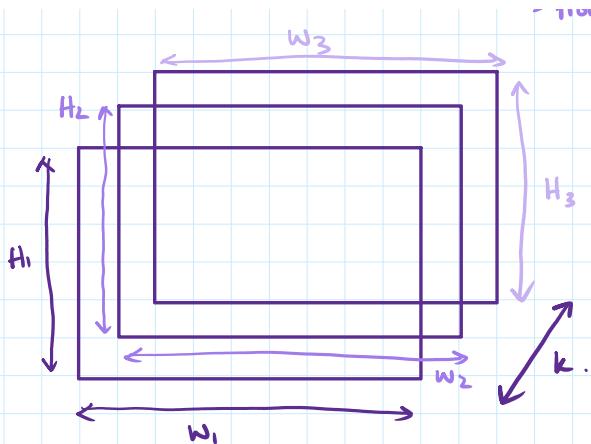
Depth = no. of filters in a layer

(filter is small matrix that slides over 1/P) 

→ How it works:

- ① Input image w/ 3 channels (R, G, B)
- ② In convolutional layer, apply multiple filters, each filter = own feature map





- (1) Input image w/ 3 channels (RGB)
- (2) In convolutional layer, apply multiple features, each filter = own feature map
- (3) If you use 32 filters, O/P feature map = 32 features.

* Size of filter = $F \times F$ → spatial extent.
no. of filters = k .
 \hookrightarrow Depth = k

Output size is just like before: where you start.

$$H_2 = \frac{W_1 + 2p - F}{s} + 1$$

$$H_2 = \frac{H_1 + 2p - F}{s} + 1$$

Now solve koshan.

Ex1: Dimension of $Y_P = 5 \times 5 \times 3$
No. of filter = 2 — k
Spatial extent = 3 — F
Zero padding = 1 — p
Stride = 2 — s

- a) Output size=?
b) Can we change stride to 3?

$$a) H_2 = \frac{W_1 + 2p - F}{s} + 1 = \frac{(5) + (2)(1) - (3)}{2} + 1 = 3$$

$$H_2 = \frac{H_1 + 2p - F}{s} + 1 = \frac{(5) + (2)(1) - (3)}{2} = 3$$

$$D_2 = k = 2$$

$$\Rightarrow O/P \text{ size} = (3 \times 3 \times 2)$$

$$b) \text{If } s=3, \quad W_2 = H_2 = \frac{4}{3} \text{ — not possible.}$$

Ex2. Input size = $288 \times 288 \times 3$
No zero padding — $p=0$
Stride = 2 — $s=2$
Output vol = $143 \times 143 \times k$
 $W_2 \quad H_2$

- a) filter size? — $F=2$
b) max pooling layer connected to o/p volume.
filter = 3×3
stride = 2
o/p volume of max pooling?
 $\hookrightarrow (W_3, H_3, D)$

$$a) W_2 = \frac{W_1 + 2p - F}{s} + 1 \quad 143 = \frac{288 - F}{2} + 1$$

$$F = 288 - 2(142) \\ \boxed{F=4}$$

$$b) \text{Input: } (143, 143, k)$$

$$s=2 \\ f=3$$

$$H_3 = \frac{H_2 + 2p - F}{s} + 1 = \frac{143 - 3}{2} + 1 = 71$$

$$W_3 = 71 \text{ too.}$$

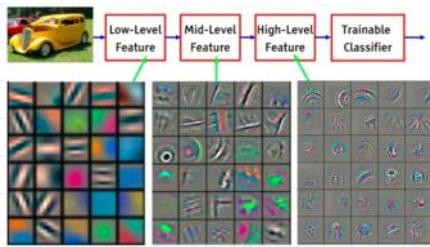
$$\Rightarrow O/P \text{ volume} = (71 \times 71 \times k)$$

Stacking layers. Why?

CNN learns more complex features

initially — simple features → build onto more high level.





initially - simple features → build onto more high level.

Also: images shrink after each layer: $\text{size}(\text{output}) < \text{size}(\text{input})$

filters @ each layer remain the same
(they find larger informed patterns)

I/P is weighted sum from prev. layers.

How many parameters in max pool layer?

none.

It doesn't need to learn anything.

$$\text{size of O/P image } O = \frac{I - ps + 1}{s}$$

Input pool size — just replace filter size (F) w/ pool size | same formula.

Calculating Number of Parameters.

W_c = no. of weights in convol. layer

B_c = no. of biases

P_c = no. of parameters.

k = size of kernels (filters)

N = no. of kernels (filters)

C = no. of channels in I/P.

$$P_c = W_c + B_c$$

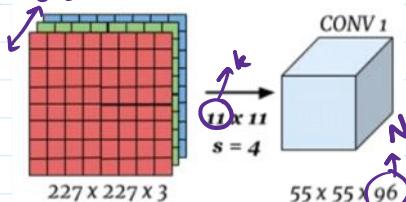
$$W_c = k^2 \times C \times N$$

each filter has $(k \times k)$ elements

$$B_c = N$$

no. of biases = no. of kernels
(cuz every filter has its own bias)

Example: $c=3$ Input:



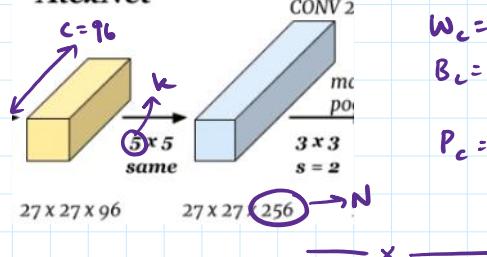
$$W_c = (11 \times 11) \times (3) \times (96) = 34348$$

$$B_c = N = 96$$

$$P_c = 96 + 34348 = 34444.$$

Example:

AlexNet



$$W_c = k^2 \times C \times N = (5 \times 5) \times (96) \times (256) = 614400$$

$$B_c = N = 256$$

$$P_c = 256 + 614400 = 614656.$$

Ok bro.

But what is a "convolution"?

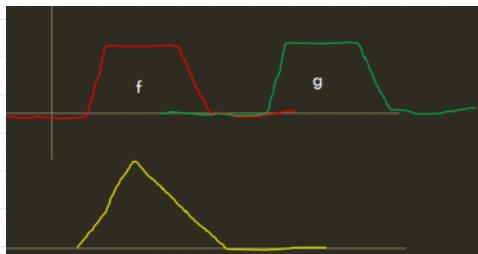
*sigh.

Conv = mechanism to blend two functions of time

$$(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau$$

"Summation of pointwise product of two functions' values subject to traversal."





summation of pointwise product of two functions values
subject to traversal!"

Now,

cross correlation = comp. of two diff time series
to detect correln b/w metrics
w/ same max & min.

auto correlation = same but w/ itself.

CNN is mostly just cross correlation.

— x —

"measure of similarity b/w a signal +
a time delayed version of another."

CNN for text:

Word embedding + CNN = text classification.

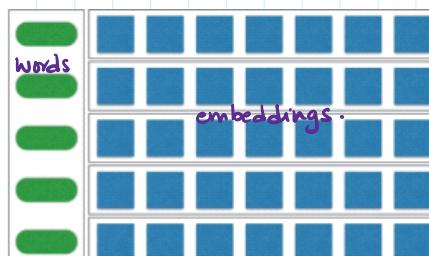
use to represent
a document.

find features
to discriminate
documents

non linearity = heck yeah

3 key aspects :

- word embedding layer
- convolution model.
- Fully connected layer (for interpretation of extracted features) — softmax.

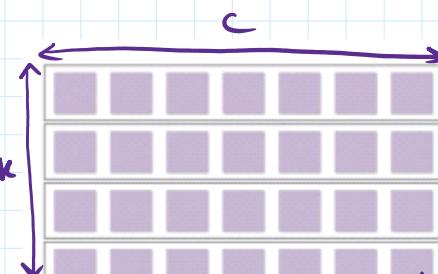


Stage 1: The representation of your Input

① Embed.

Multiple filters? Can.
Same length? Can.

produce multiple
filter maps
that's all.



Stage 2: This is a filter of length 4

② Can make 2 also.

K = no. of rows

C = no. of cols.

→ This is embedding b/w

CNN - 1D vs 2D Convolution

	1D Convolution	2D Convolution
Input Data	Sequential data (e.g., time series, audio signals)	Image data (e.g., RGB images, grayscale images)
Input Shape	(Length, Channels)	(Height, Width, Channels)
Filter Shape	(Filter Length, Channels)	(Filter Height, Filter Width, Channels)
Operation	Slides filters over the sequence	Slides filters over the 2D spatial dimensions of the image
Stride	Step size for moving the filter along the sequence	Step size for moving the filter along height and width
Padding	Adds values (usually zeros) at the start and end of sequence	Adds values (usually zeros) around the image border
Output Shape	(Output Length, Output Channels)	(Output Height, Output Width, Output Channels)
Computational Complexity	Generally less complex due to fewer dimensions involved	Generally more complex due to higher dimensional data
Applications	Time series analysis, speech recognition, text processing	Image classification, object detection, image segmentation