

Week #4 :- Special applications: Face recognition & Neural Style Transfer

Face Recognition:-

What is face recognition?

Face detection + Liveliness \Rightarrow Face recognition.

Liveliness: whether the picture is real human ^{himself} or not.

Face Verification vs Face Recognition:

Verification:-

1:1 problem.

- Input: image or name/ID
- Output: Whether the input image is that of claimed person

Recognition:-

- Has a database of K persons
- Get an Input: Image
- Output: ID if the image is any of the K persons
(or "not recognized")

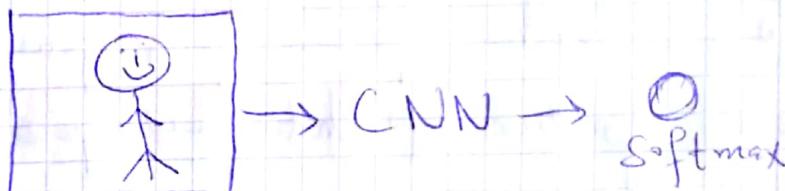
Let's see how can we build a face verification system and we will use that to develop face recognition.

One Shot Learning :- What this means:-

"You need to be able to recognize an image given just one single image"

→ actually what we have learned so far is that more the training examples, more the accuracy of your algorithm will be.

One approach :-



This doesn't work because, what if the amount of the persons in your team is not fixed what if another person joins your team, do you have to train your Convnet again? Seems like, not a good idea.

Instead You would do something like:
Learning a "similarity" function

$d(\text{img1}, \text{img2}) = \text{degree of difference b/w images}$

So, if two images are of same person's you want it to output a small number otherwise a very large number.

→ Threshold.

if $d(\text{img1}, \text{img2}) \leq T$ "Same person"

" " " > T "Different person"

That's how you address a verification problem

To use this for recognition task :-

Given the input image, you compare that with
with the images from your data base using the
 $d(\text{img}^1, \text{img}^2)$ function that will be learned
and that's how the Recognition problem is solved.

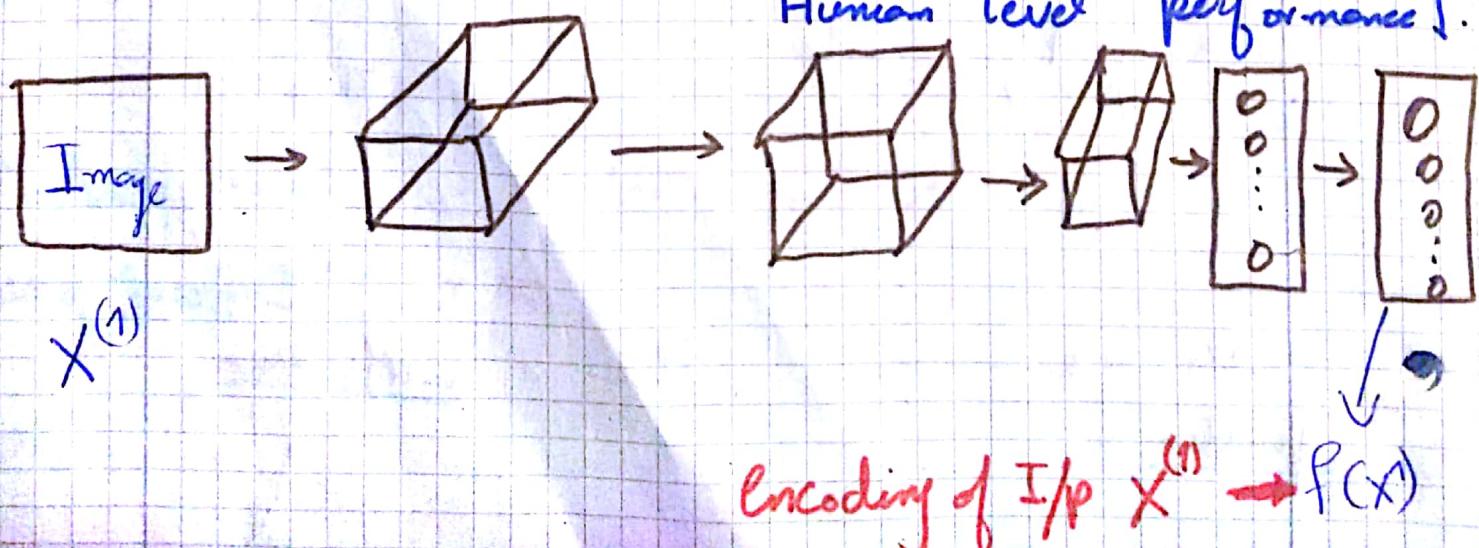
- Now if another person joins your team, you just add
another picture of that new person and it just
works fine.
- Let's see how you can actually train a neural
network to learn this function "d"

Siamese Network:-

Job of the 'd' function is to input two images
and tell you how close and far they are from
being similar.

Good way to do it is through Siamese
network.

[Taigman et. al., 2014. DeepFace closing the gap to
Human level performance].



The way you build a face recognition is the way you want to recall compare two pictures:-

What you could do, is feed the second picture $x^{(2)}$ with some parameters to some neural network and get a different encoding of that picture called $f(x^{(2)})$.

$$d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_2^2$$

Now how do you train this Siamese network?

So, you want to train a N.N such that the mapping or encoding it computes results in the function d .

Parameters of NN define an encoding $f(x^i)$

Learn parameters so that:

if $x^{(i)}, x^{(j)}$ are the same person,

$\|f(x^{(i)}) - f(x^{(j)})\|_2^2$ is small.

if $x^{(i)}, x^{(j)}$ are the different persons,

$\|f(x^{(i)}) - f(x^{(j)})\|_2^2$ is large.

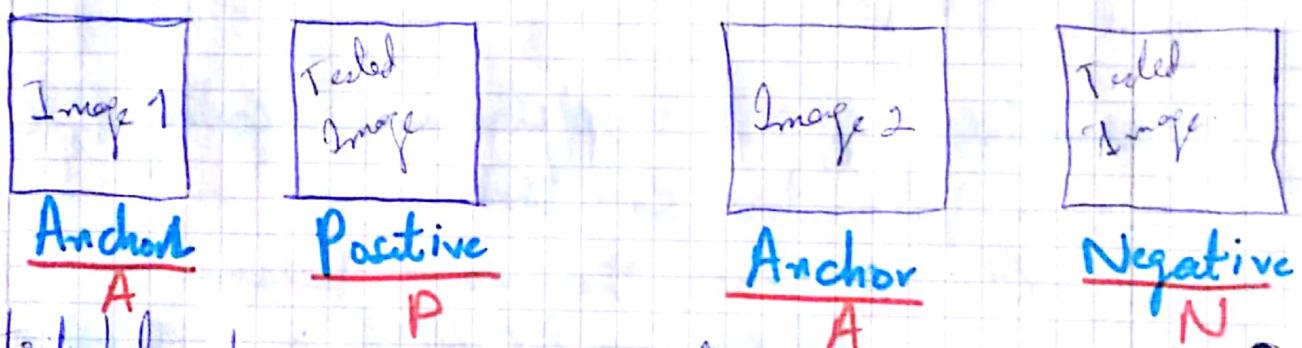
Now, how do you define an objective function for this?

Triplet Loss: One good way to learn the parameters of the neural network is to define the objective function called "Triplet loss" and apply gradient descent on that.

[Schroff et al., 2015, FaceNet: A unified embedding for face recognition and Clustering]

To apply the triplet loss, you need to compare pairs of images.

e.g. Given a picture, to learn the parameters of a neural network, you have to look at several pictures at the same time.



In triplet loss, you are always going to look at one anchor image and you want the distance b/w anchor and positive image (meaning same person) and the distance of anchor from the negative image (meaning different person) much further apart.

This is what gives rise the term triplet loss.
You always will be looking at 3 images, anchor, positive and negative image.

Your work:

$$\frac{\|f(A) - f(p)\|^2}{d(A, p)} \leq \frac{\|f(A) - f(N)\|^2}{d(A, N)}$$

$$\|f(A) - f(p)\|^2 - \|f(A) - f(N)\|^2 \leq 0 \quad \text{--- (1)}$$

- One trivial way to satisfy this is to learn everything $f(A), f(p), f(N) = 0$.
- Another trivial way to satisfy the expression is to have identical values like $f(A) = f(p) = f(N)$.

So, in order to prevent from above 2 cases happening, we need to make sure that not only the expression from the left hand side is ≤ 0 but it should be quite a bit smaller than that such that

$$\|f(A) - f(p)\|^2 - \|f(A) - f(N)\|^2 \leq 0 \quad \text{--- (2)}$$

another hyperparameter

This prevents the network from outputting trivial solutions.

Normally γ is on left hand side.

$$\|f(A) - f(p)\|^2 - \|f(A) - f(N)\|^2 + \gamma \leq 0 \quad \text{--- (2)}$$

Margin.

So, e.g. if $d(A, p) = 0.5$ and

$d(A, N) = 0.51$ then without margin the

(1) equation would be satisfied, but ~~without~~ that's not a good criteria we want $d(A, N)$ to be 0.7 at least or maybe something much bigger.

So, to ensure that, we set some margin that our algorithm needs to overcome in order to give valid that person is same.

→ Let's take the equation (2) and formalize it to define the triplet loss.

Loss function:-

Given 3 images $\boxed{A, P, N}$

$$L(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \epsilon)$$

You want the blue underlined expression to be less than zero, if it is then we will get 0 in output and the output is negative, otherwise (the output is positive if that is not greater) less than or equal to zero.

Cost function:-

$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)})$$

Suppose Training set: 10K pictures of 1K persons.

What you will do is, take 10K pictures and use it to generate / select triplets and train your algorithm by applying gradient descent.

You need pairs of (A, P) , so for this algo. training you need multiple pictures atleast for some persons.

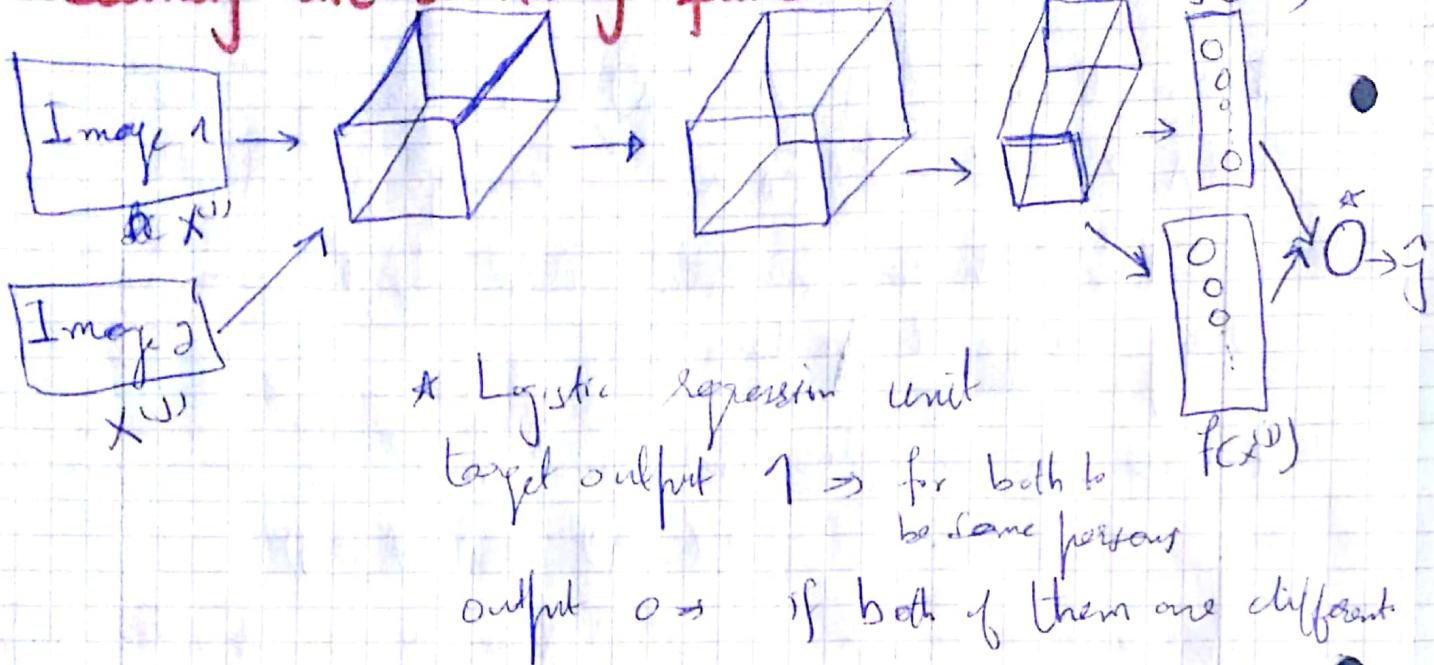
Choosing the triplets A, P, N

- During training if A, P, N are chosen randomly,
 $d(A, P) + \alpha \leq d(A, N)$ is easily satisfied.
because chances are that , A and N can be
very different pictures.
Choose the triplets that're "hard" to train on
so that $d(A, P) \approx d(A, N)$.
- This way algorithm will have to try harder to
try to keep the thing on the left down and
the right one up. So that, there is at least
a mapping of left or right side.
- Choosing random triplets won't let gradient
descend any work and you ② identifiy
will be easily satisfied. Choosing it carefully
will push the gradient descent to work
harder

Face Verification and Binary Classification:-

Besides using a triplet loss, there is another way to perform this task. Let's see how face recognition can be easily posed as binary classification problem.

Learning the similarity function:-



\rightarrow face recognition as binary classification problem \Rightarrow alternative to triplet loss.

Now what does first logistic regression do?

$$\hat{y} = \sigma \left(\sum_{k=1}^{\text{total units in FC}} w_k \underbrace{\left| f(x^{(1)})_k - f(x^{(2)})_k \right| + b}_{\text{weight}} \right)$$

Compute and see whether the images are of same or different persons.

\rightarrow Other variations of blue underlined expression:-

$$\frac{(f(x^{(1)})_k - f(x^{(2)})_k)^2}{f(x^{(1)})_k + f(x^{(2)})_k} \times 2$$

\rightarrow Kullback-Leibler encoding.

So:-

Input : Pair of images ; Output 0/1. whether same or different).

One computational ~~task~~^{hack} is that:

- You compute embeddings of new person's image and the one already in your database and compare them.
- You can pre-compute the embeddings of the images already in your database and then you won't have to compute it again and again.
- Face Verification: Supervised Learning.

Input: pairs of images

Output: 0/1 \Rightarrow use backprop for learning.

Neural Style Transfer:

What is Neural Style Transfer?

Take one image and another image of an animation or artistic style. What neural style transfer will do is that, it will allow you to generate an image which will be the 1st original image painted in style of the second one.

Let's define Notations

Image to Style \Rightarrow Content (C)

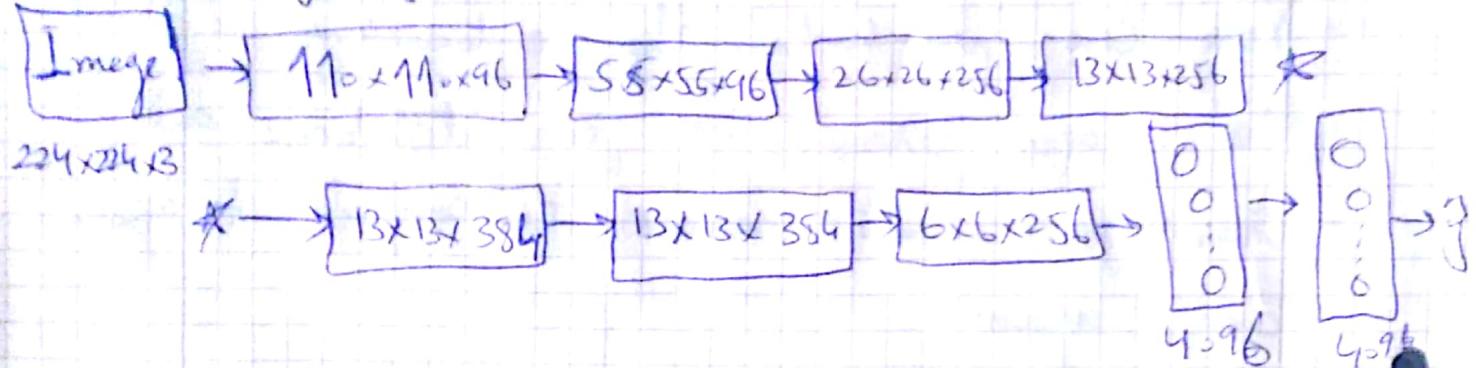
Style Image \Rightarrow Style (S)

Generated Image \Rightarrow (G)

So, in order to do that, let's see what are deep ConvNets actually doing?

What are deep ConvNets doing?

Let's say you have created a ConvNet



→ We want to visualize what the hidden layers in this network are computing.

→ Let's start like this:-

Pick a unit in layer 1. Find the 9 image patches that maximize the unit's activation.

In other words:-

Pass your training set through the ConvNet and figure out, which is the image that maximizes that particular unit's activation.

You'll find small patches of 9 images (under consideration) that will affect the maximization of that unit.

→ Now you can pick a different hidden unit in the layer and do the same thing.

→ You'll see different patterns observed by different hidden

units. One main thing you can observe is that, longer
it is detecting simple things such as edges or
particular shade of color.

[Zeiler and Fergus., 2013, Visualizing and understanding Convolutional networks]

What if you do this for some of the deeper layers
hidden in the network.

- In deeper layer, a hidden unit will see a
larger region of the image

So as a result units in the deeper layers will
start detecting more complex patterns.

Now that we have somewhat proper understanding of how
does CNNs work, let's use this intuition in
order to build full algorithm.

Cost function :-

To build a neural style transfer system let's
~~also~~ define a cost function for the generated
image and later you'll see that, by reducing
that cost function we can generate the image
we want.

$$\text{Content}(C) + \text{Style}(S) \rightarrow \text{Generated Image. } (G)$$

Apply $J(G)$ to estimate how good is a
generated image and use gradient descent
to optimize it.

We will define 2 part Cost function :-

$$J(G) = \alpha J_{\text{Content}}(C, G) + \beta J_{\text{Style}}(S, G)$$

α Checks the similarity b/w Content of generated image and style image

β Check similarity b/w Content of generated image and style image

$\alpha, \beta \Rightarrow$ weights to specify the relative weighting b/w content cost and style cost.

[Gatys et al., 2015], A neural algorithm of artistic style, Images on slides generated by Justin Johnson]

Find the generated image G

1. Initiate G randomly

$G: 100 \times 100 \times 3 \Rightarrow$ any size

2. Use gradient descent to minimize $J(G)$

$$G := G - \frac{\partial}{\partial G} J(G).$$

→ Randomly initialized image will be a random noise and with every iteration it will converge.

Now let's define Cost function.

Content Cost function:-

- Overall Cost function is :-

$$J(G) = \alpha J_{\text{Content}}(C, G) + \beta J_{\text{Style}}(S, G)$$

- Say you use hidden layer l , to compute content cost.
- If layer l is very earlier layer then large content of your generated image would be very similar to your content image.
- If l is very deep layer then, it's just asking, "Well if there is a dog in your content image then make sure there is a dog somewhere in your generated image." So, in practice layer l is chosen somewhere in between.
- Use pre-trained ConvNet (e.g. VGG Network)
- Let $a^{[l]}(C)$ and $a^{[l]}(G)$ be the activation of layer l on the images
- If $a^{[l]}(C)$ and $a^{[l]}(G)$ are similar, both images have similar content

$$J_{\text{Content}}(C, G) = \frac{1}{2} \| a^{[l]}(C) - a^{[l]}(G) \|_F^2$$

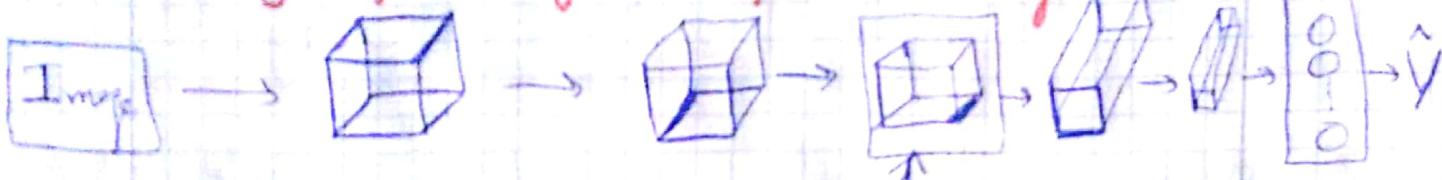
Normalization
Constant

Element wise.
Sum of squares
of differences

- Shows how similar are content and generated images in terms of content.

Style Cost function :-

: Meaning of "Style" of an image:-



- Say you are using layer l's activation to measure "style"

What we are going to do is to:

→ Define Style as Correlation b/w activation across channels in layer l.

Here is what I mean?

Suppose you like the layer's activations and we are gonna ask how correlated are the activations across different channels.

To understand this phrase let's

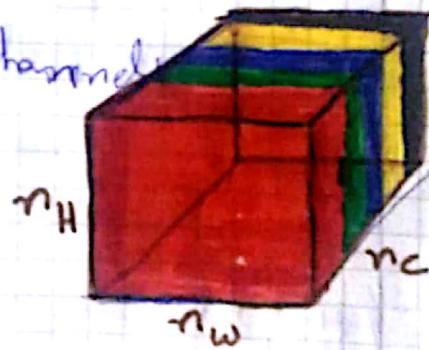
look at the channels of this

block. For the sake of

understanding let's say we have

5 block's layers

→ Take red channel and green channel and say how correlated are the activations in first 2 channels.



What you do is look at values at different positions of these two channels and see how much correlation these two channels show.

Now, Why does this capture style?

Recall, we were trying to get an intuition about how ConvNets compute what they compute.

Let's say red channel corresponds to the neuron we are trying to figure out if there is a little vertical texture in the particular position.

and second (green) channel corresponds to the neuron trying to figure out orangish color.

Now, what does this mean for this channels to be highly correlated?

This means that, whatever part of the image will have vertical texture, it will also have orangish color. and opposite is the case for them to be uncorrelated.

→ So, Correlation tells you that, whichever of these high level texture components tend to occur or not occur together and degree of correlation tells you, how often these high level features occur together or not in different parts of the images.

If degree of correlation tells you the similarity then applying the same concept you can access how much style in generated image is correlated with the style of the style image at certain place of an image.

Let's formalize intuition :-

Given image we compute **Style Matrix** which will measure all those correlation we have discussed.

Let $\{a\}$

$a_{ijk}^{[l]} = \text{activation at } (i, j, k)$. $G^{[l]}$ is $n_c \times n_c$

$n_c \Rightarrow$ channels

So $G_{kk'}^{[l]}$ \Rightarrow Compares the correlation b/w activation of layers l b/w channel k and k' .

$K = 1, \dots, n_c^{[l]}$

Un-normalized

Cross Co-variance

$$G_{kk'}^{[l](S)} = \sum_i \sum_j a_{ijk}^{[l](S)} a_{ijk'}^{[l](S)}$$

S \Rightarrow style

If activations \Rightarrow both of them are large then $G_{kk'}^{[l]}$ will be large, if they are uncorrelated then it might be small.

$G \Rightarrow$ generated image

$$G_{kk'}^{[l](G)} = \sum_{i=1}^{n_h} \sum_{j=1}^{n_w} a_{ijk}^{[l](G)} a_{ijk'}^{[l](G)}$$

Cost function

$$J_{\text{style}}^{[l]}(S, G) = \|G^{[l](S)} - G^{[l](G)}\|_F^2$$

$$J_{\text{style}}^{[l]}(s, g) = \frac{1}{(n_h^{[l]} n_w^{[l]} n_c^{[l]})^2} \sum_k \sum_{k'} (G_{kk'}^{[l]} - G_{kk'}^{[l](g)})^2$$

→ Finally, it turns out that, you get more pleasing results if you use the style cost function from multiple different layers. So overall style cost function can be defined as:

$$J_{\text{style}}(s, g) = \sum_l \lambda^{[l]} J_{\text{style}}^{[l]}(s, g)$$

So overall cost function:

$$J(g) = \alpha J_{\text{content}}(c, g) + \beta J_{\text{style}}(s, g)$$

Use gradient descent or other sophisticated algorithm on it to minimize the cost.

1D and 3D Generalizations:-

Just like we do convolutions in 2D images

like $(16 \times 16 \times 3) * [5 \times 5 \times 2] \Rightarrow [16 \text{ filters}] \rightarrow [1 \times 1 \times 16] * [5 \times 5 \times 6] \rightarrow [6 \times 6]$

1D data

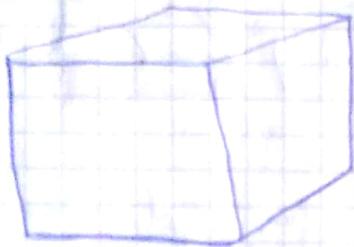
$$\begin{bmatrix} 1 & 2 & 0 & 1 & 5 & 3 & 1 & 8 & 1 & 2 & 4 & 1 & 7 \end{bmatrix} *$$



$$16 \times 1 * [5 \times 1] \xrightarrow{6 \text{ filters}} [1 \times 16] * [5 \times 16] \xrightarrow{32 \text{ filters}} [6 \times 32]$$

Mostly for 1D data Recurrent NN are used.

How about a 3d data?



3D Volume

$$[14 \times 14 \times 16]$$

One channel.



3D g. 1ke

$$[5 \times 5 \times 5]$$

$$\rightarrow [10 \times 10 \times 10]$$

e.g. 3D data
from CT
Scan.

This filter detects feature from 3D data.