

Machine learning

Lect 11

Dr. Haran

Mukherjee



SVM Numerical

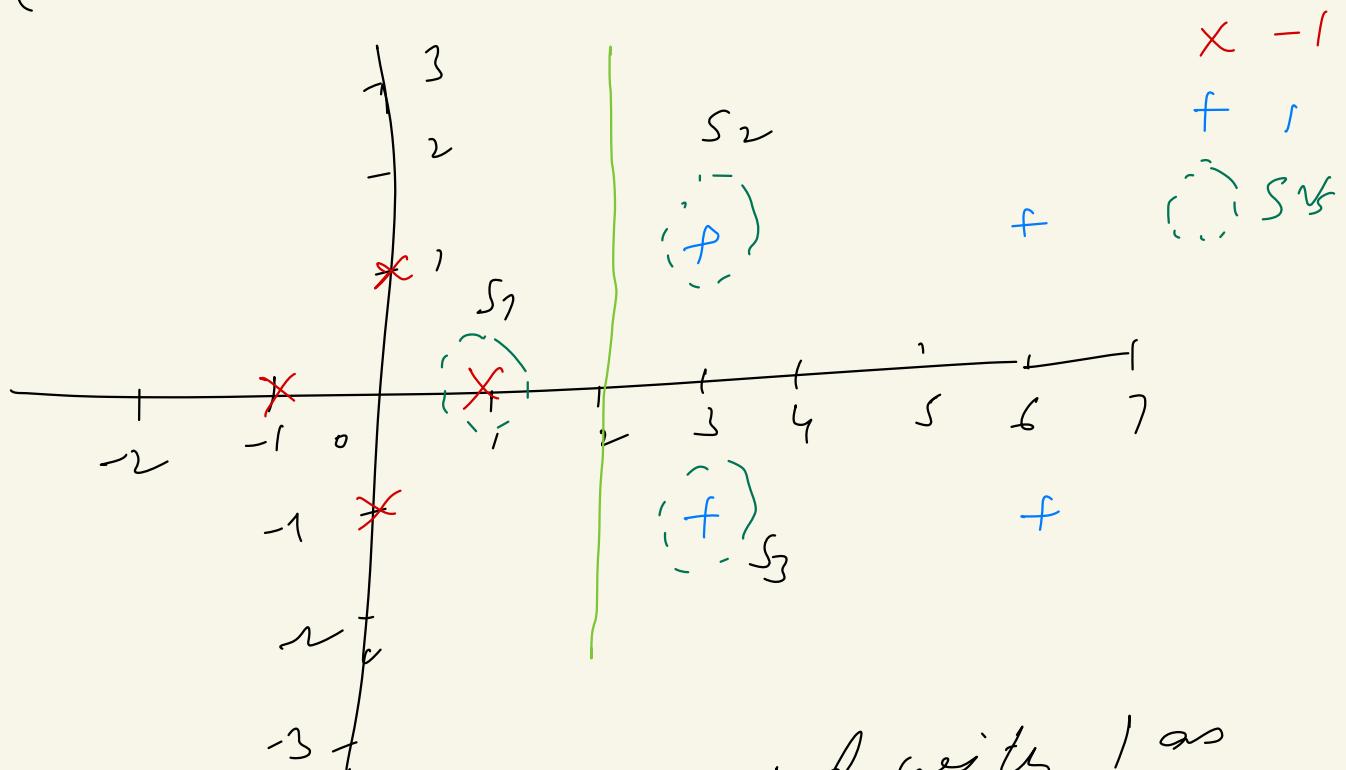
Suppose we are given the following data points

$$\left\{ \begin{pmatrix} 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 3 \\ -1 \end{pmatrix}, \begin{pmatrix} 6 \\ 0 \end{pmatrix}, \begin{pmatrix} 6 \\ -1 \end{pmatrix} \right\}$$

five labeled +1

$$\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix} \right\}$$

unlabeled -1



Each vector is augmented with 1 as bias input

$$S_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \tilde{S}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$$

$$S_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix} \quad \tilde{S}_2 = \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \quad S_3 = \begin{pmatrix} 3 \\ -1 \end{pmatrix} \quad \tilde{S}_3 = \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix}$$

$$\mathcal{L}_1 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} = -1$$

$$\cancel{\mathcal{L}_1} \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} = +1$$

$$\mathcal{L}_1 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} \tilde{s}_1 \\ \tilde{s}_2 \\ \tilde{s}_3 \end{pmatrix} = +1$$

$$\mathcal{L}_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = -1$$

$$\mathcal{L}_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} = +1$$

$$\mathcal{L}_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \mathcal{L}_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \mathcal{L}_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} = +1$$

$$\begin{aligned} \alpha_1(1+0+1) + \alpha_2(3+0+1) \\ + \alpha_3(3+0+1) = -1 \end{aligned}$$

$$\begin{aligned} \alpha_1(3+0+1) + \alpha_2(9+1+1) \\ + \alpha_3(9-1+1) = 1 \end{aligned}$$

$$\begin{aligned} \alpha_1(3+0+1) + \alpha_2(9-1+1) + \\ \alpha_3(9+1+1) = 1 \end{aligned}$$

$$2\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1$$

$$4\alpha_1 + 11\alpha_2 + 9\alpha_3 = 1$$

$$4\alpha_1 + 9\alpha_2 + 11\alpha_3 = 1$$

$$\begin{aligned} \alpha_1 = -3.5 & \quad \alpha_2 = 0.75 \\ \alpha_3 = 0.75 & \end{aligned}$$

$$\begin{aligned} \tilde{\omega} &= \sum_i \alpha_i s_i \\ &= -3.5 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \end{aligned}$$

$$= \begin{pmatrix} 1 \\ 0 \\ -2 \end{pmatrix}] \rightarrow \tilde{\omega}$$

sinc

$$y = \omega x + b$$

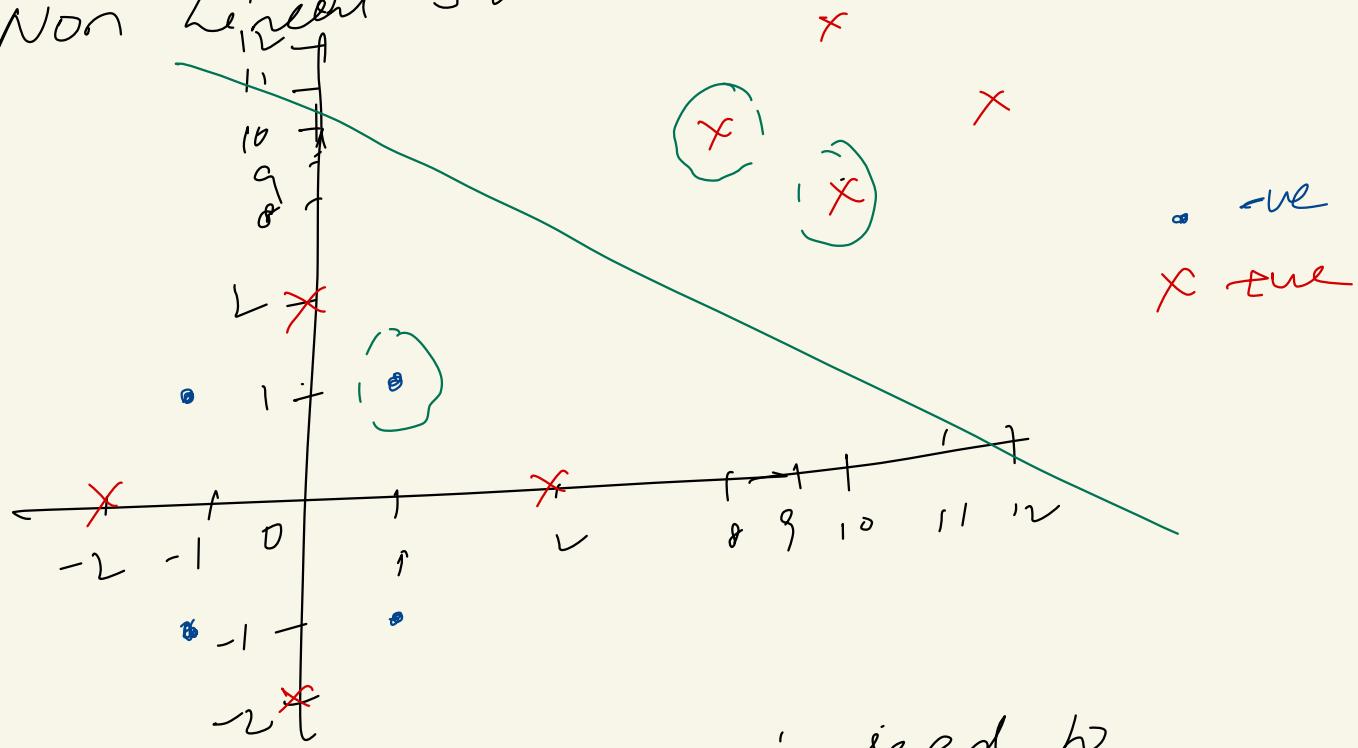
$$\omega = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad b = -2$$

R -ve class points $\rightarrow \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ -1 \end{pmatrix}$

the +ve class points $\rightarrow \begin{pmatrix} 4 \\ 0 \end{pmatrix}, \begin{pmatrix} 5 \\ 1 \end{pmatrix}, \begin{pmatrix} 5 \\ -1 \end{pmatrix}, \begin{pmatrix} 6 \\ 0 \end{pmatrix}$

find the generality function

Non Linear SVM



$\phi \rightarrow$ Mapping function is reqd to transform these data to a new feature space where a separating hyperplane can be found.

$$\phi(n_1, n_2) = \begin{cases} \begin{pmatrix} 6 - n_1 + (n_1 - n_2)^2 \\ 6 - n_2 + (n_1 - n_2)^2 \end{pmatrix} & \text{if } n_1 + n_2 \geq 0 \\ \begin{pmatrix} n_1 \\ n_2 \end{pmatrix} & \text{otherwise} \end{cases}$$

\rightarrow Blue class vectors are: $\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix},$

$\begin{pmatrix} -1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ no change since

$$\sqrt{n_1^2 + n_2^2} < 2 \text{ ff vectors}$$

→ Red class vectors $\begin{pmatrix} 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} -2 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ -2 \end{pmatrix}$

$$\phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \phi \begin{pmatrix} 2 \\ 0 \end{pmatrix} = \begin{pmatrix} 6 - 2 + (2-0)^2 \\ 6 - 0 + (2-0)^2 \end{pmatrix} = \begin{pmatrix} 8 \\ 10 \end{pmatrix}$$

$$\phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \phi \begin{pmatrix} 0 \\ 2 \end{pmatrix} = \begin{pmatrix} 10 \\ 8 \end{pmatrix}$$

$$\phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \phi \begin{pmatrix} -2 \\ 0 \end{pmatrix} = \begin{pmatrix} 12 \\ 10 \end{pmatrix}$$

$$\phi \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \phi \begin{pmatrix} 0 \\ -2 \end{pmatrix} = \begin{pmatrix} 10 \\ 12 \end{pmatrix}$$

$$s_1 = \begin{pmatrix} 8 \\ 10 \end{pmatrix} \quad s_2 = \begin{pmatrix} 10 \\ 8 \end{pmatrix} \quad s_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\tilde{s}_1 = \begin{pmatrix} 8 \\ 10 \\ 1 \end{pmatrix} \quad \tilde{s}_2 = \begin{pmatrix} 10 \\ 8 \\ 1 \end{pmatrix} \quad \tilde{s}_3 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$d_1 \tilde{s}_1 \tilde{s}_1 + d_2 \tilde{s}_2 \tilde{s}_1 + d_3 \tilde{s}_3 \tilde{s}_1 \underset{\sim}{\sim} \underset{\sim}{\sim} = +1$$

$$d_1 \tilde{s}_2 \tilde{s}_2 + d_2 \tilde{s}_2 \tilde{s}_2 + d_3 \tilde{s}_2 \tilde{s}_2 \underset{\sim}{\sim} \underset{\sim}{\sim} = +1$$

$$d_1 \tilde{s}_3 \tilde{s}_1 + d_2 \tilde{s}_3 \tilde{s}_2 + d_3 \tilde{s}_3 \tilde{s}_3 \underset{\sim}{\sim} \underset{\sim}{\sim} = -1$$

$$165\omega_1 + 161\omega_2 + 19\omega_3 = +1$$

$$164\omega_1 + 165\omega_2 + 19\omega_3 = +1$$

$$19\omega_1 + 19\omega_2 + 3\omega_3 = -1$$

$$\omega_1 = \omega_2 = 0.859 \quad \omega_3 = -1.4219$$

$$\tilde{\omega} = \sum_i \omega_i s_i$$

$$\tilde{\omega} = \omega_1 \begin{pmatrix} 8 \\ 10 \\ 1 \end{pmatrix} + \omega_2 \begin{pmatrix} 10 \\ 8 \\ 1 \end{pmatrix} + \omega_3 \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} 0.1243 \\ 0.1243 \\ -1.2501 \end{pmatrix}$$

$$y = \omega^n + b$$

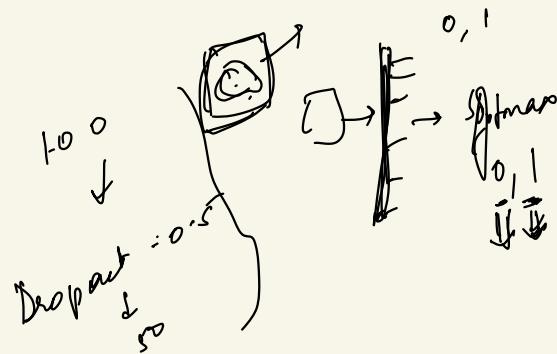
$$\omega = \begin{pmatrix} 0.1243 & 0.1243 \\ 0.1243 & 0.1243 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$b = -1.2501 / 0.1243 \\ = -10.057$$

Test point $n_1, n_2 = -1, 2$ belongs to red point class five

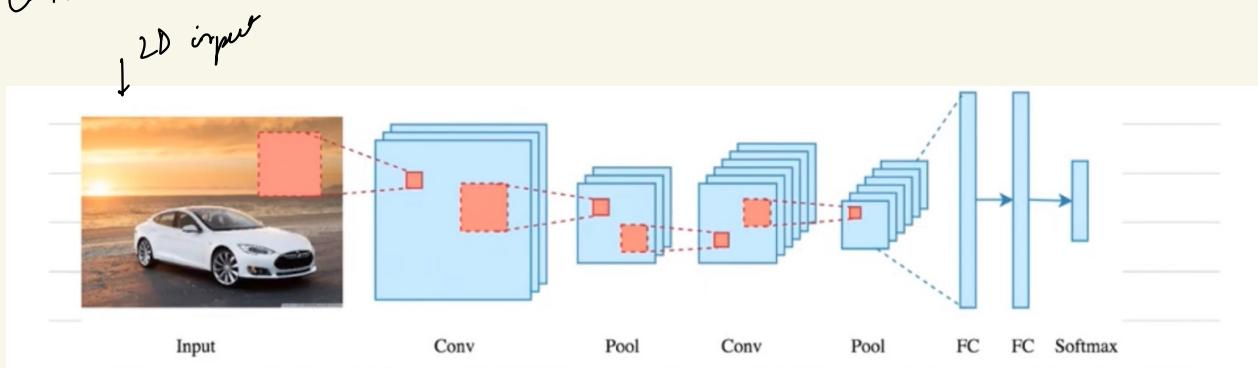
$$\phi(n_1, n_2) = \begin{pmatrix} 16 \\ 13 \end{pmatrix}$$

$$w \phi(n) = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \begin{pmatrix} 16 \\ 13 \end{pmatrix} = 29710$$



Introduction to Deep Learning (CNNs)

Convolutional Neural Network



Correlation
Convolution
Input
1D
input
 $h(x)$
f θ^0

CNNs over ANNs - i) No context learnt (Neighborhood of pixels)
ii) More parameters to learn in ANNs

no flipping \downarrow flip kernel 18°
Correlation / Convolution
In case of CNNs \rightarrow it is actually correlation only
Conv (filter size, stride, padding) \rightarrow pooling

① Filter size

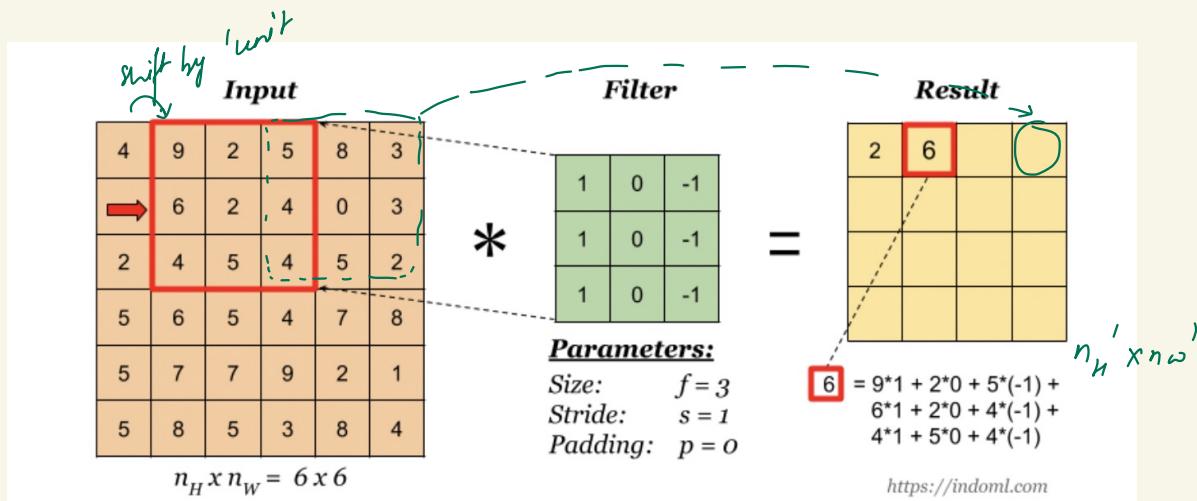
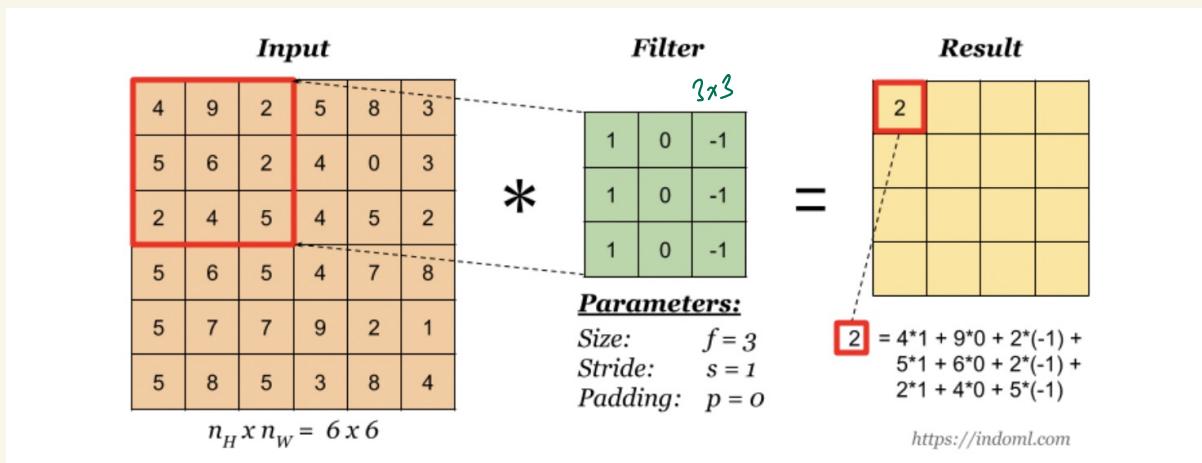
1	0	1
0	0	0
-1	0	-1

3×3

2	3	4	2	1
6	9	2	1	6
2	8	7	2	7
1	1	2	6	1
3	3	5	3	2

8×8

② Convolution



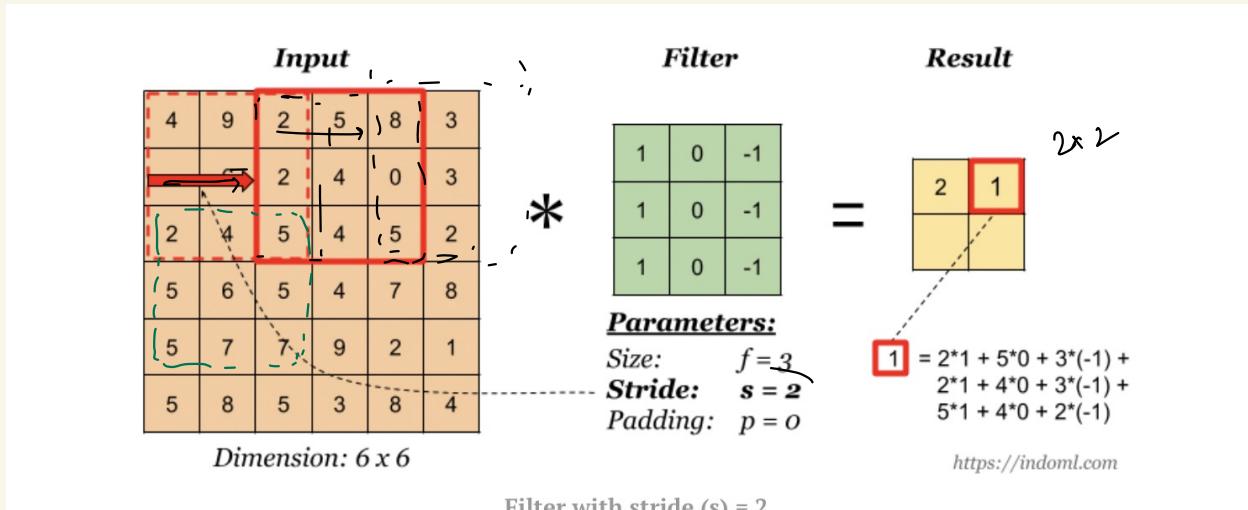
$$n_H' = n_H - f + 1 = \\ = 6 - 3 + 1 = 4$$

stride

$$w = 5-1, p=0$$

$$n_H' = \frac{n_H - f + 1}{s} + 1$$

$$n_H' = \left\lfloor \frac{n_H + 2p - f + 1}{s} \right\rfloor$$



Filter with stride (s) = 2

Padding

$$\begin{matrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix} \quad 4 \times 4 \rightarrow 6 \times 6 \rightarrow 8 \times 8$$

$p = 1 \quad p = 2$

since $p=1$ increase image size by 2^{units}

$$n_H' = n_H - f + 2p + 1$$

Convolution reduces size, padding = valid \downarrow same $\Rightarrow p = \frac{f-1}{2}$
 reduces size of image \downarrow keeping original size of image same $\Rightarrow p = \frac{f-1}{2} - 1$
 $6 \times 6 \rightarrow 6 \times 6$

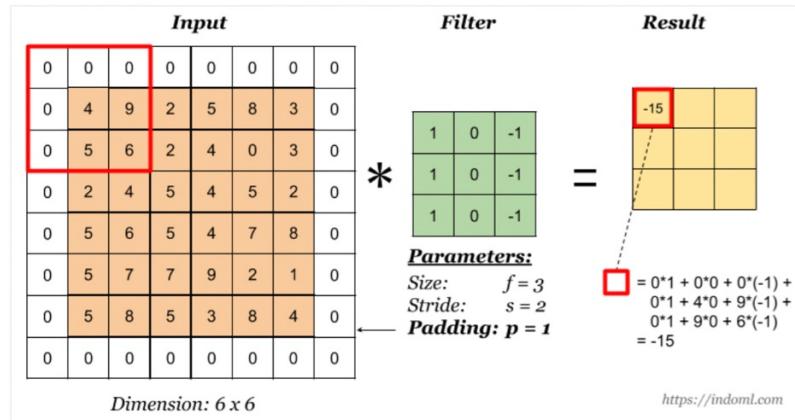
$$n_H' = \left\lfloor \frac{n_H + 2p - f}{s} + 1 \right\rfloor = \underline{\underline{n_H}}$$

$$= R = \frac{f-1}{2}$$

Padding

Padding has the following benefits:

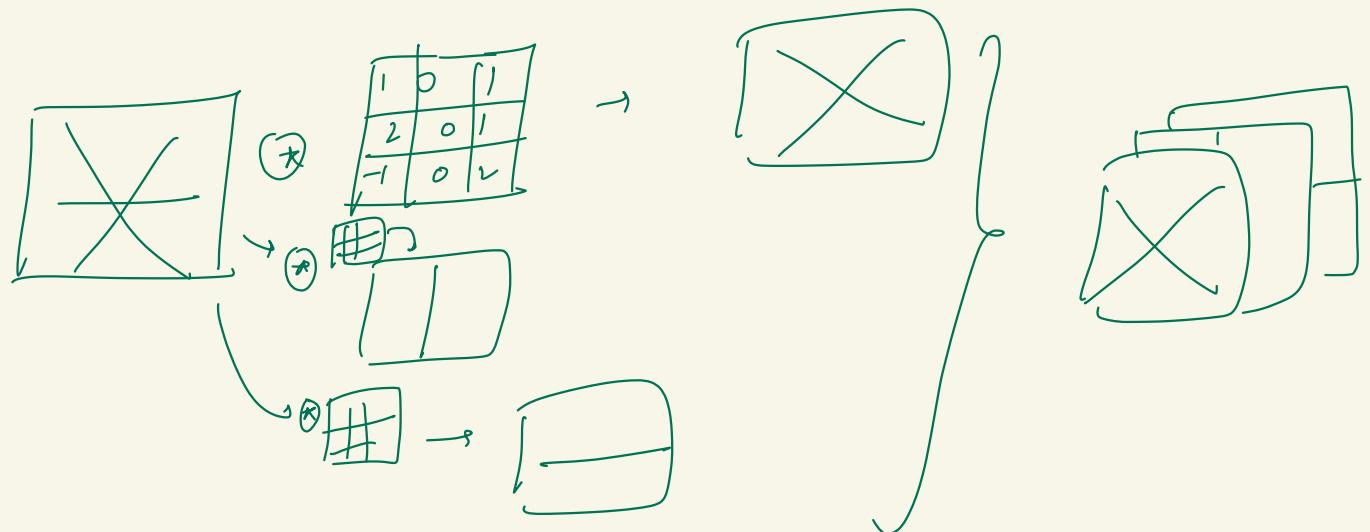
1. It allows us to use a CONV layer without necessarily shrinking the height and width of the volumes. This is important for building deeper networks, since otherwise the height/width would shrink as we go to deeper layers.
2. It helps us keep more of the information at the border of an image. Without padding, very few values at the next layer would be affected by pixels as the edges of an image.

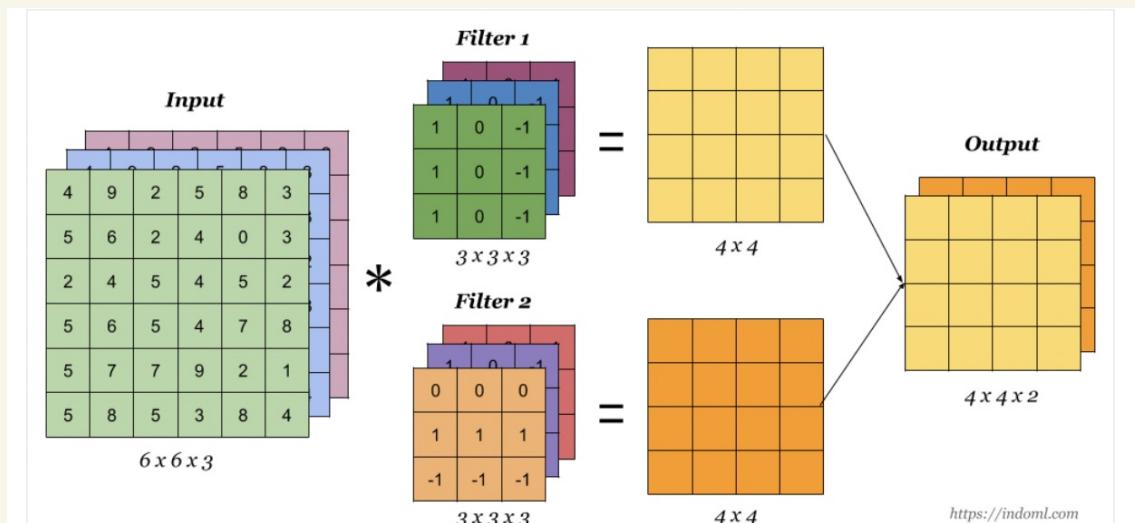
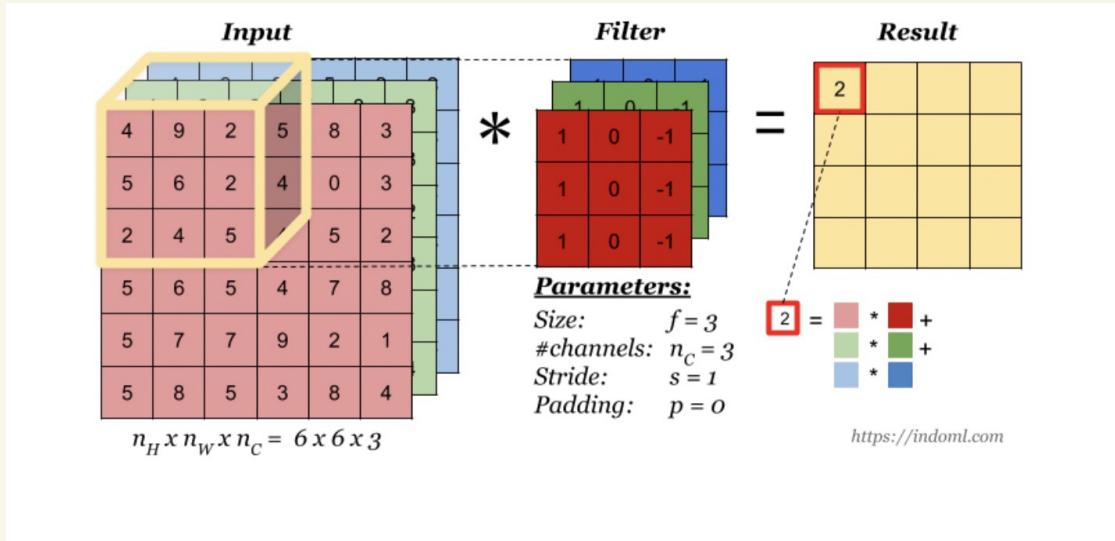


Notice the the dimension of the result has changed due to padding. See the following section on how to calculate output dimension.

Some padding terminologies:

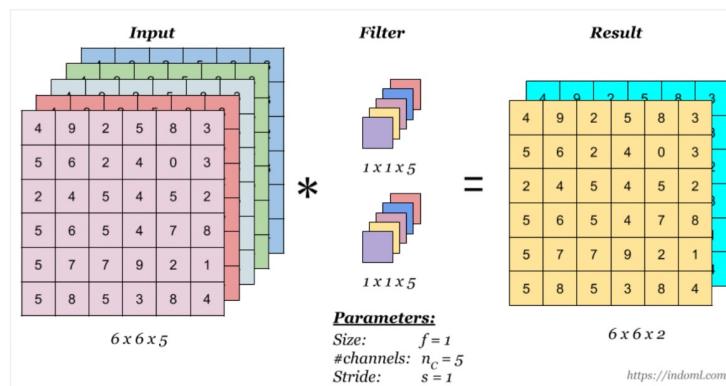
- “**valid**” padding: no padding
- “**same**” padding: padding so that the output dimension is the same as the input

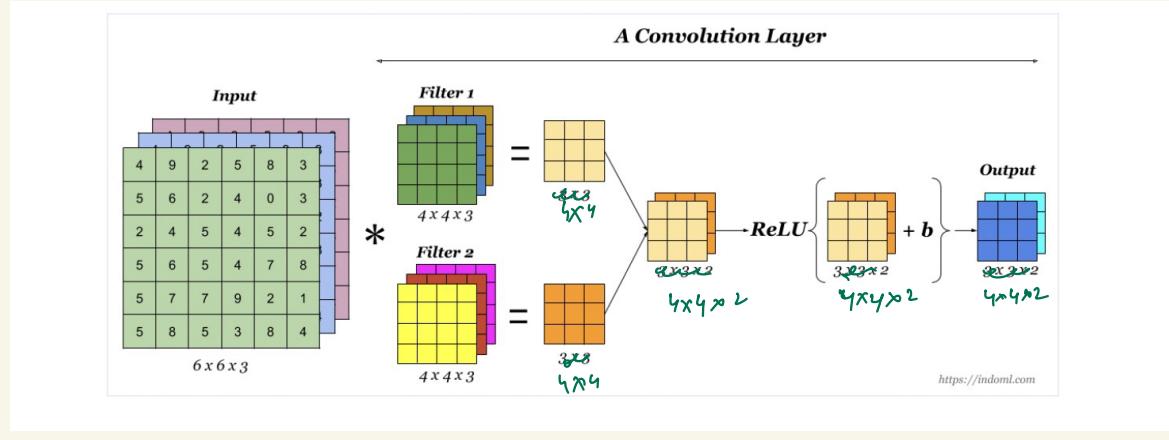




1 x 1 Convolution

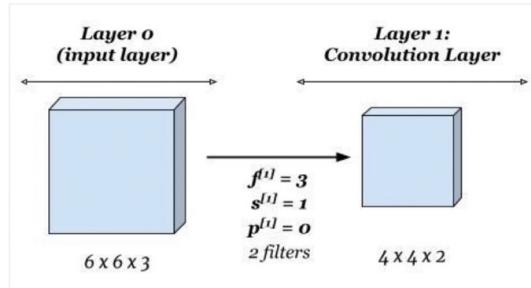
This is convolution with 1×1 filter. The effect is to flatten or “merge” channels together, which can save computations later in the network:





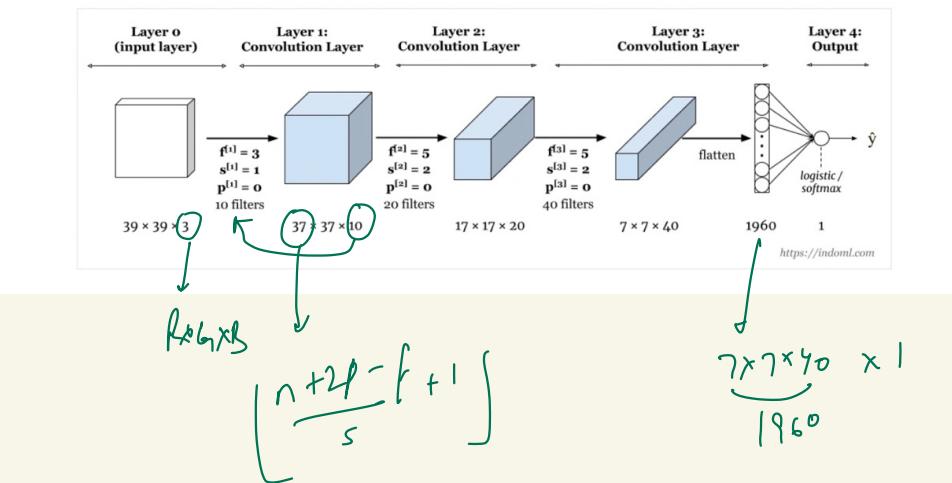
Shorthand Representation

This simpler representation will be used from now on to represent one convolutional layer:

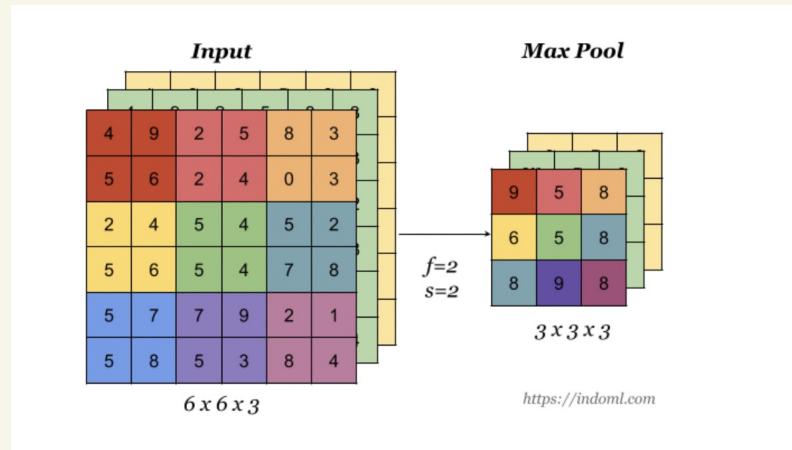
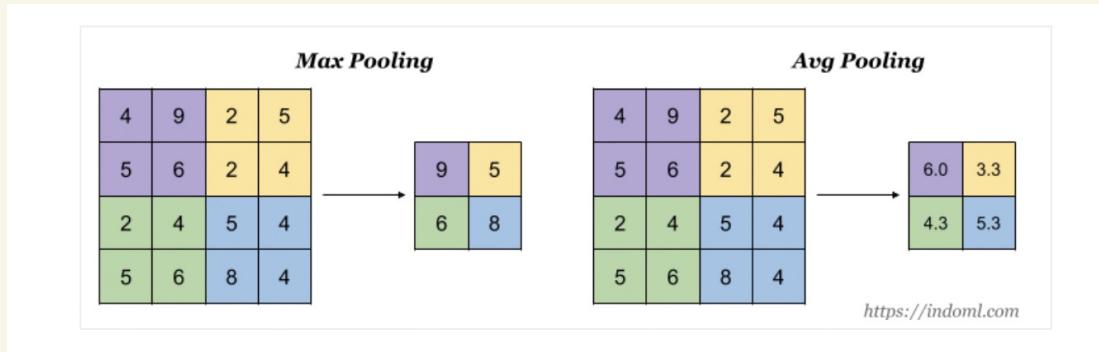


Sample Complete Network

This is a sample network with three convolution layers. At the end of the network, the output of the convolution layer is flattened and is connected to a logistic regression or a softmax output layer.



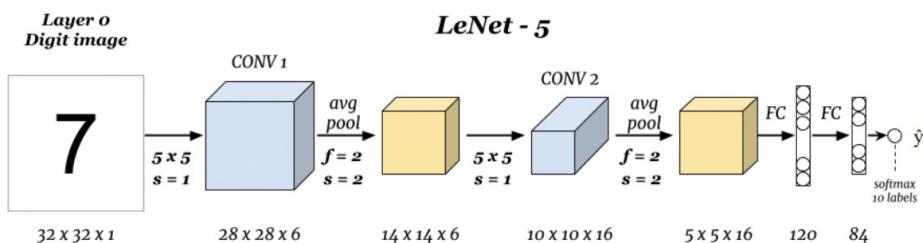
Pooling
 ↗ Avg
 ↗ Max (more used)
 My parameters: size (f) , stride(s) , type (max / avg)



Well Known Architectures

Classic Network: LeNet – 5

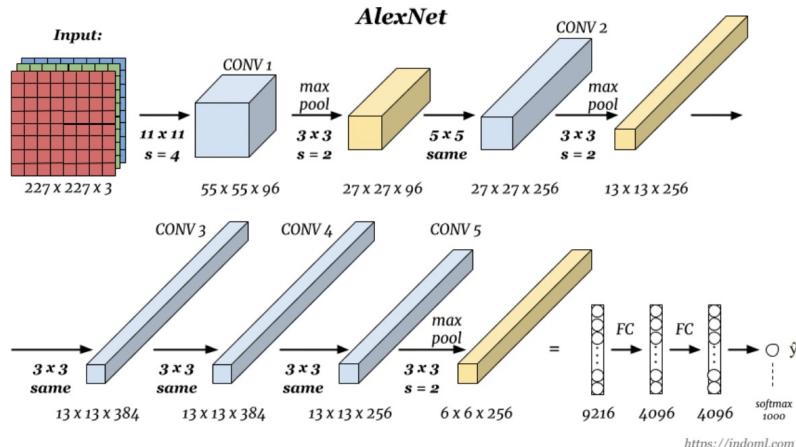
One example of classic networks is LeNet-5, from *Gradient-Based Learning Applied to Document Recognition* paper by Y. Lecun, L. Bottou, Y. Bengio and P. Haffner (1998):



- Number of parameters: ~ 60 thousands.

Classic Network: AlexNet

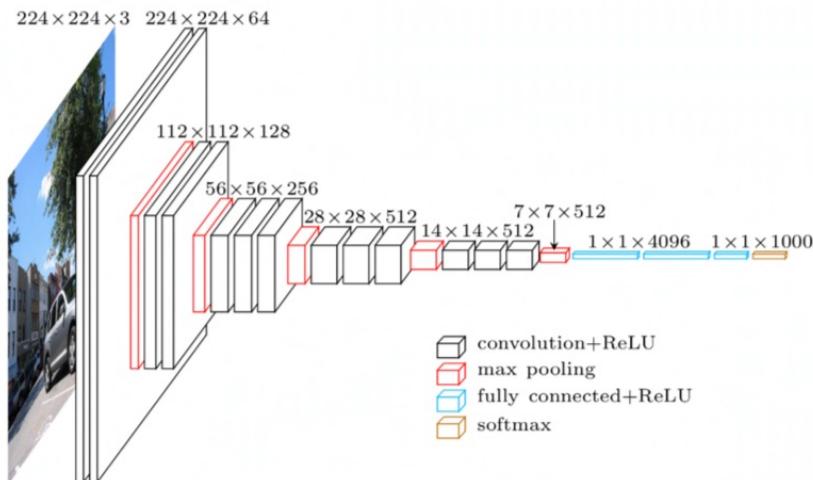
AlexNet is another classic CNN architecture from *ImageNet Classification with Deep Convolutional Neural Networks* paper by Alex Krizhevsky, Geoffrey Hinton, and Ilya Sutskever (2012).



- Number of parameters: ~ 60 millions

Classic Network: VGG-16

VGG-16 from *Very Deep Convolutional Networks for Large-Scale Image Recognition* paper by Karen Simonyan and Andrew Zisserman (2014). The number 16 refers to the fact that the network has 16 trainable layers (i.e. layers that have weights).

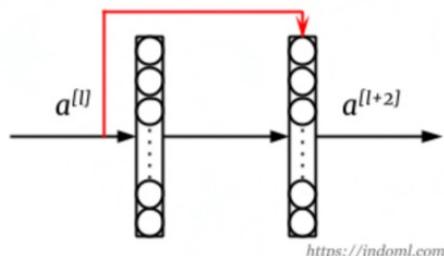


(image from blog.heuritech.com)

- Number of parameters: ~ 138 millions.
- The strength is in the simplicity: the dimension is halved and the depth is increased on every step (or stack of layers)

ResNet

The problem with deeper neural networks are they are harder to train and once the number of layers reach certain number, the training error starts to raise again. Deep networks are also harder to train due to exploding and vanishing gradients problem. ResNet (Residual Network), proposed by He et al in [Deep Residual Learning for Image Recognition paper](#) (2015), solves these problems by implementing skip connection where output from one layer is fed to layer deeper in the network:



In the image above, the skip connection is depicted by the red line. The activation $a^{[l+2]}$ is then calculated as:

$$z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]}$$

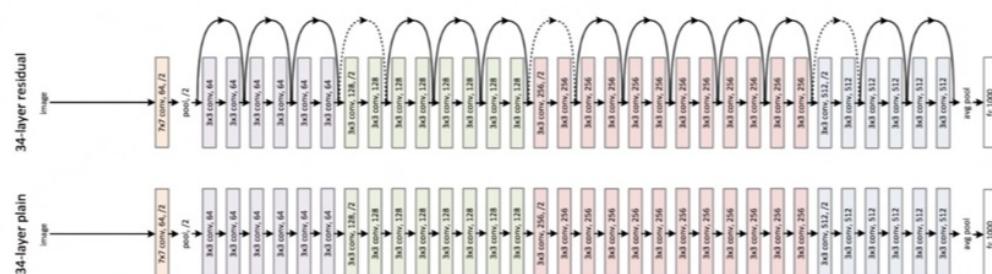
$$a^{[l+2]} = g^{[l+2]}(z^{[l+2]} + \color{red}{a^{[l]}})$$

The advantages of ResNets are:

- performance doesn't degrade with very deep network
- cheaper to compute
- ability to train very very deep network

ResNet works because:

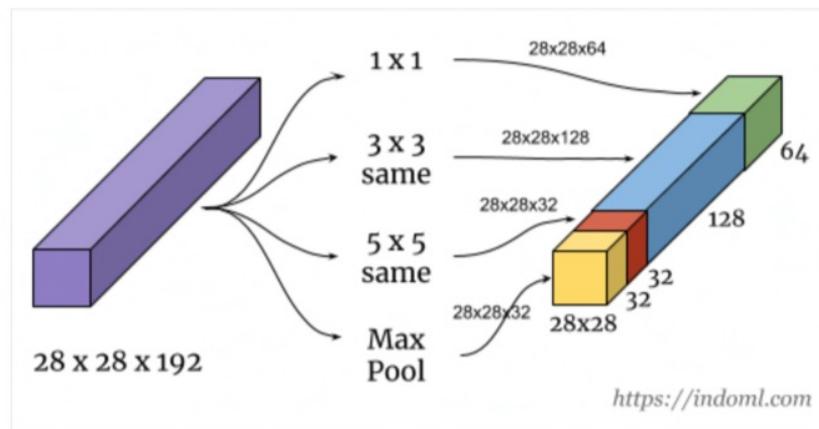
- identity function is easy for residual block to learn
- using a skip-connection helps the gradient to back-propagate and thus helps you to train deeper networks



Comparison of 34 layers ResNet with plain network (image from euler.stat.yale.edu)

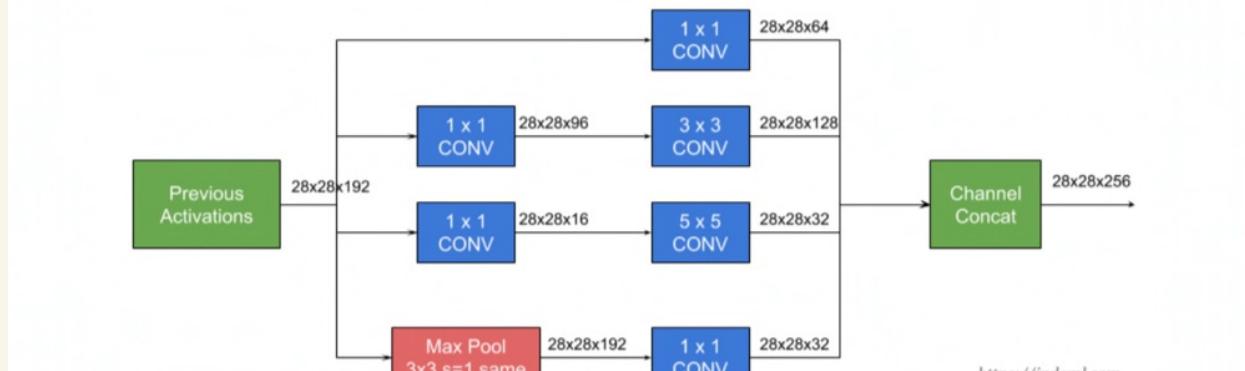
Inception

The motivation of the inception network is, rather than requiring us to pick the filter size manually, let the network decide what is best to put in a layer. We give it choices and hopefully it will pick up what is best to use in that layer:

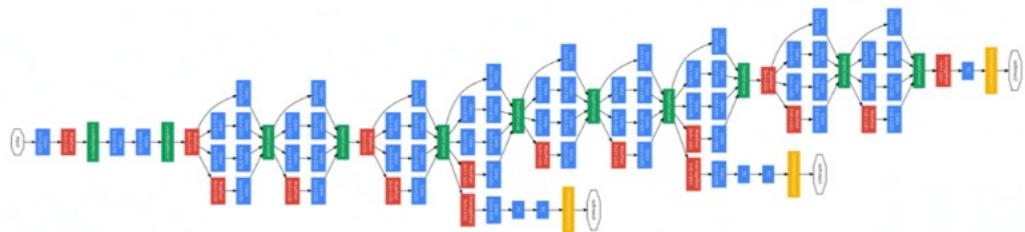


The problem with the above network is computation cost (e.g. for the 5×5 filter only, the computation cost is $(28 \times 28 \times 32) \times (5 \times 5 \times 192) = \sim 120$ millions).

Using 1×1 convolution will reduce the computation to about 1/10 of that. With this idea, an inception module will look like this:



Below is an inception network called **GoogLeNet**, described in *Going Deeper with Convolutions paper* by Szegedy et all (2014), which has 9 inception modules:



GoogLeNet architecture (image [source](#))

Data Augmentation

- Mirroring
- Random cropping
- Less common (perhaps due to their complexity):
 - Rotation
 - Shearing
 - Local warping
- Color shifting \rightarrow 4

