

# Computer Vision

COURSE 4  
WEEK 1

1

Image Classification  
Object detection  
Neural Style Transfer

} Computer  
vision  
problems

filter = kernel

# Edge Detection

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

Vertical filter

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

horizontal filter

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

Sobel filter

3	0	-3
10	0	-10
3	0	-3

Scharr filter

Convolutional neural network  
frameworks

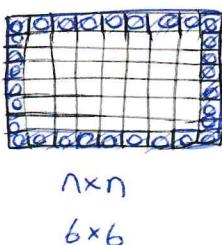
python: conv-forward

tensorflow: tf.nn.conv2d

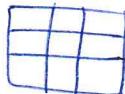
Keras : conv2D

$$n \times n * f \times f = n-f+1 \times n-f+1$$

# Padding



\*  
↓  
asterisk



$$= n+2p-f+1 \times n+2p-f+1 \quad p = \text{padding}$$

$6 \times 6$

Problems was ; without padding

- Shrinking output (dimensions)
- throwing away a lot of info from the edges of image

# Valid and Same Convolutions:

Valid: no padding  $n \times n * f \times f \rightarrow n-f+1 \times n-f+1$

Same: pad so that output size is the same as the input size,

$$n+2p-f+1 \times n+2p-f+1$$

$$n+2p-f+1 = n$$

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

$f$  is usually odd.

$3 \times 3$  very common    $5 \times 5$     $7 \times 7$



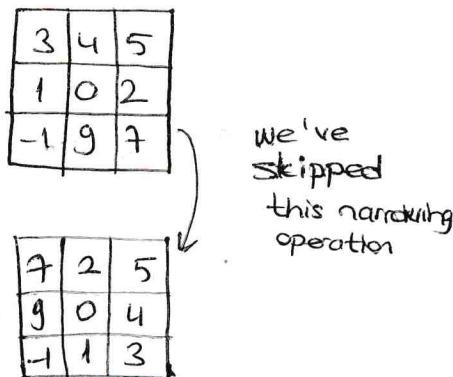
## Strided Convolutions

COURSE 4  
WEEK 1

2

$$n \times n * f \times f \longrightarrow \frac{n+2p-f}{s} + 1 \times \frac{n+2p-f}{s} + 1$$

Technical note on cross-correlation vs convolution



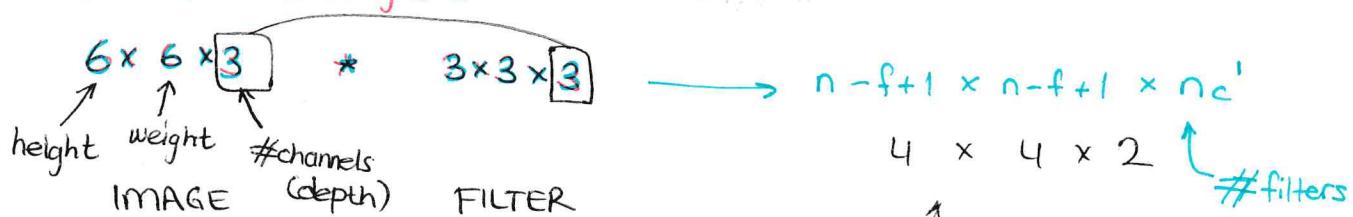
In deep learning literature it is called as a convolutional operation.

Technically this operation is maybe better called cross-validation.

$$(A * B) * C = A * (B * C) \rightarrow \text{"associativity" in maths}$$

## Convolutions Over Volume

Convolutions on RGB images :



Multiple Filters :

$$\begin{aligned} 6 \times 6 \times 3 * 3 \times 3 \times 3 &= \begin{array}{c} \text{vertical edge filter} \\ \text{horizontal edge filter} \end{array} \\ * 3 \times 3 \times 3 &= \begin{array}{c} 2D \\ 4 \times 4 \\ 4 \times 4 \end{array} \longrightarrow 4 \times 4 \times 2 \end{aligned}$$

## One Layer Of a Convolutional Network Example

$$\begin{aligned} 6 \times 6 \times 3 * 3 \times 3 \times 3 &\longrightarrow \text{ReLU } ((4 \times 4) + b_1) \rightarrow 4 \times 4 \\ \underbrace{a^{[0]}}_{a^{[0]}} * \underbrace{3 \times 3 \times 3}_{W^{[1]} \cdot a^{[0]}} &\longrightarrow \text{ReLU } ((4 \times 4) + b_2) \rightarrow 4 \times 4 \\ z^{[1]} = W^{[1]} a^{[0]} + b^{[1]} & \\ a^{[1]} = g(z^{[1]}) & \end{aligned}$$

$"W^{[1]} a^{[0]}"$

$4 \times 4 \times 2$   
# filters

Number of parameters in one layer :

If you have 10 filters that are  $3 \times 3 \times 3$  in one layer of neural network, how many parameters does that layer have?

$$3 \times 3 \times 3 = 27 + 1 = 28 \quad \text{bias} \quad 28 \times 10 = 280 \text{ parameters}$$

## Summary of Notation

COURSE 4  
WEEK 1

3

If layer  $l$  is a convolution layer:

$f^{(l)}$  = filter size

$P^{(l)}$  = padding

$s^{(l)}$  = stride

$n_c^{(l)}$  = number of filters

Each filter is :  $f^{(l)} \times f^{(l)} \times n_c^{(l-1)}$

Activations :  $a^{(l)} \rightarrow n_H^{(l)} \times n_W^{(l)} \times n_c^{(l)}$

Weights :  $f^{(l)} \times f^{(l)} \times n_c^{(l-1)} \times n_c^{(l)}$

bias :  $n_c^{(l)} \rightarrow (1, 1, 1, n_c^{(l)})$  #filters in layer  $l$

Input :  $n_H^{(l-1)} \times n_W^{(l-1)} \times n_c^{(l-1)}$

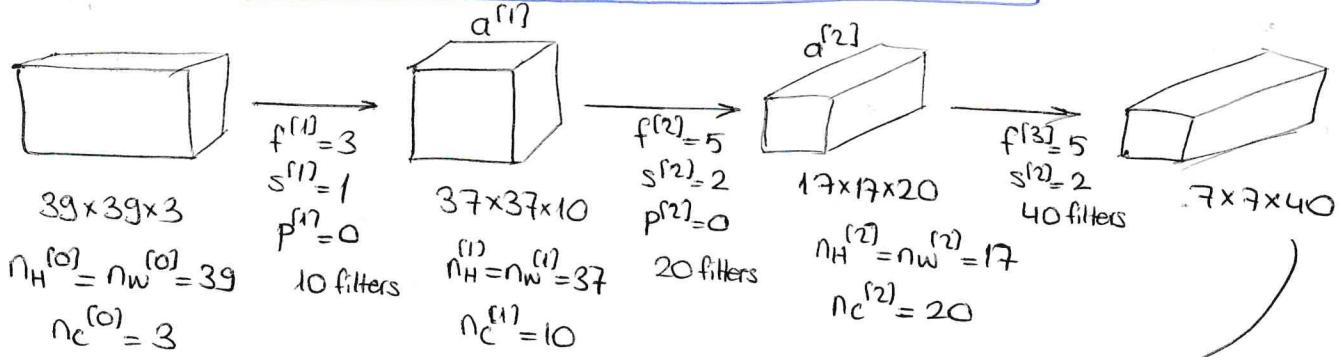
Output :  $n_H^{(l)} \times n_W^{(l)} \times n_c^{(l)}$

$$n_{HW}^{(l)} = \frac{n_{HW}^{(l-1)} + 2P^{(l)} - f^{(l)}}{s^{(l)}} + 1$$

$$A^{(l)} = m \times n_H^{(l)} \times n_W^{(l)} \times n_c^{(l)}$$

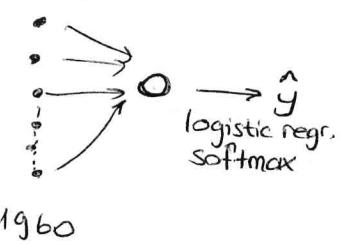
$\underbrace{\hspace{10em}}$   
 $n_c \times n_H \times n_W$

### Simple Convolutional Network Example



### Types of layer in convolutional network:

- Convolution (CONV)
  - Pooling (POOL)
  - Fully Connected (FC)
- } Simpler than CONV to define



### Pooling Layers

Pooling layer : Max pooling

$$5 \times 5 \times 2 \rightarrow 3 \times 3 \times 2 \quad \left( \frac{n+2p-f}{s} + 1 \right)$$

$\downarrow$   
#channel  
( $n \times n \times n_c$ )

$\downarrow$   
#channel  
( $3 \times 3 \times n_c$ )

### Average pooling

In very deep neural network, you might use average pooling to collapse your representation.

$$n_H \times n_W \times n_c \rightarrow \left[ \frac{n_H-f}{s} + 1 \right] \times \left[ \frac{n_W-f}{s} + 1 \right] \times n_c$$

One interesting property of max pooling is that it has a set of parameters to learn. There is actually nothing for gradient descent to learn.

\* max pooling is more common than average pooling.

### Hyperparameters:

$f$ : filter size

$s$ : stride

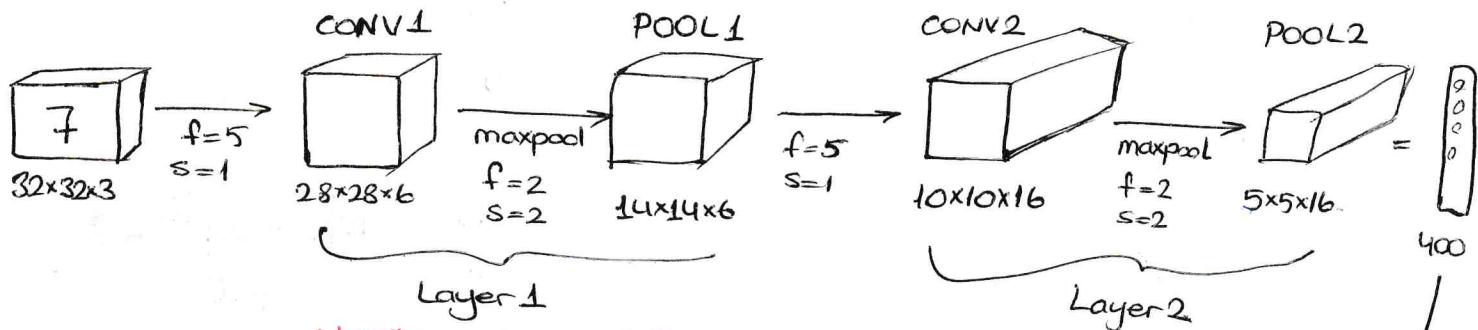
Max or average pooling

$f=2, s=2$

$f=3, s=2$

$p$ : padding (rarely)

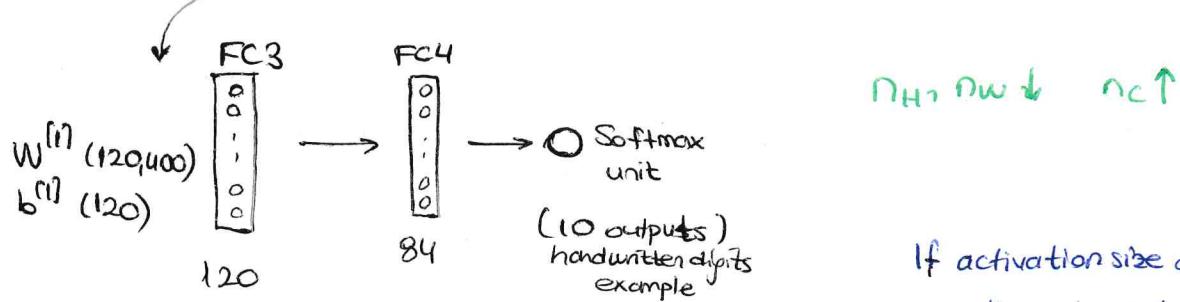
(LeNet-5) → Inspired by this net.



Usually people record the number of layers that have weight, that have parameters.

The pooling layer has no parameters (weights)

CONV1 and POOL1 come together (are grouped) as layer 1



If activation size drops too quickly, that is usually not great for performance

### Why Convolutions?

! Number of parameters in the conv layer remains quite small.

Reasons ~~of~~ how small ~~number~~ of parameters run well :

**Parameter sharing**: A feature detector (such as vertical edge detector) that is useful in one part of the image is probably useful in another part of the image.

**Sparcity of connections**: In each layer, each output value depends only on a small number of inputs.

Why look at case studies?

Good way to get intuition on how to build convnets is to read or to see other examples of effective conv nets.

Outline:

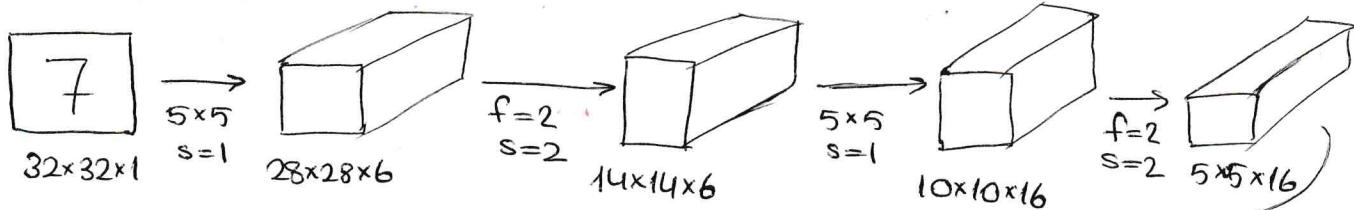
Classic networks :

- LeNet-5
- AlexNet
- VGG -16

Resnet (152)

Inception

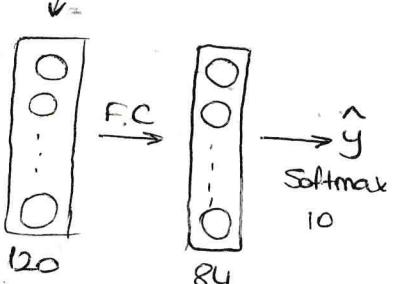
} these ideas are important to get some intuition and very interesting and very helpful for your work.

Classic NetworksLeNet-5

60K parameters :

$n_H, n_W \downarrow n_C$

conv pool      conv  
(or more)      (or more)      pool      fc      fc      output

Advanced comments:

In the paper people used sigmoid and tanh nonlinearities.  
not ReLU

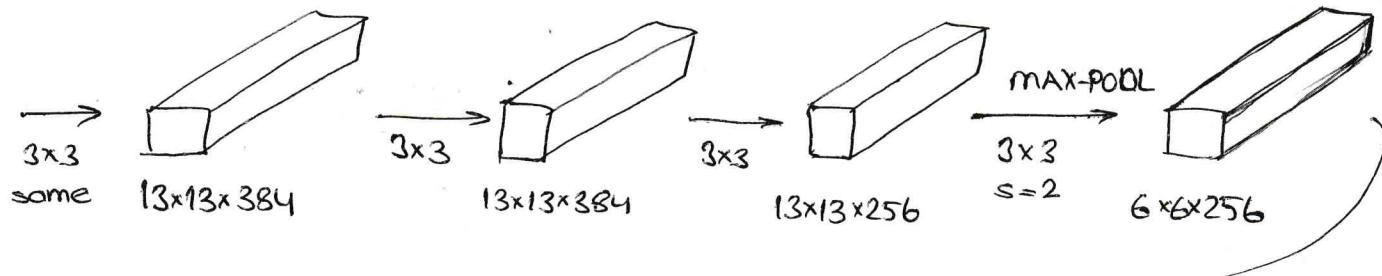
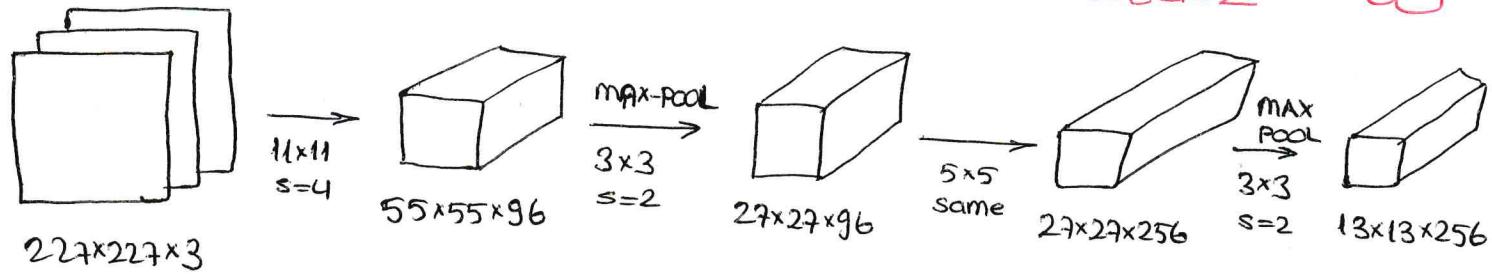
computers were slow

non-linearity after pooling layer

these ideas are from section 2 and 3 of the paper.  
about architecture

has bunch of experiments and results

### AlexNet :



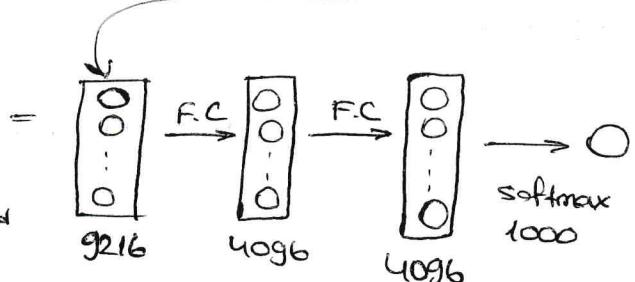
- similarity to LeNet but much bigger

$\sim 60 \text{ M parameters}$

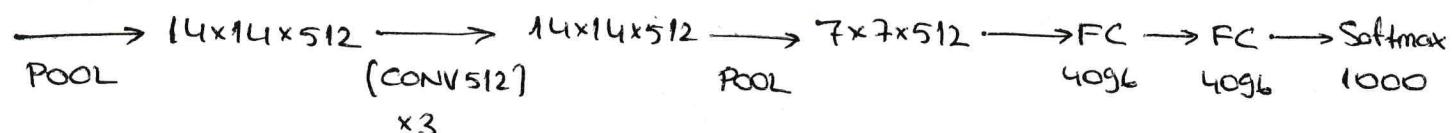
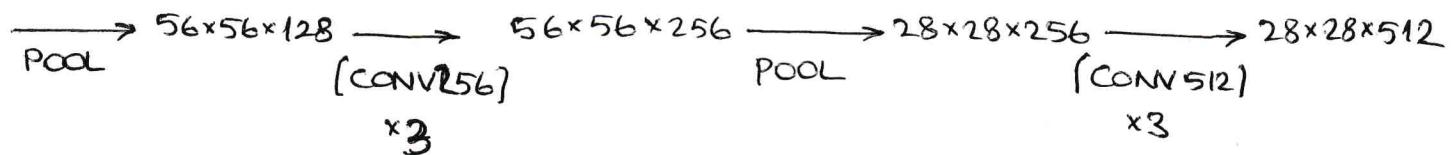
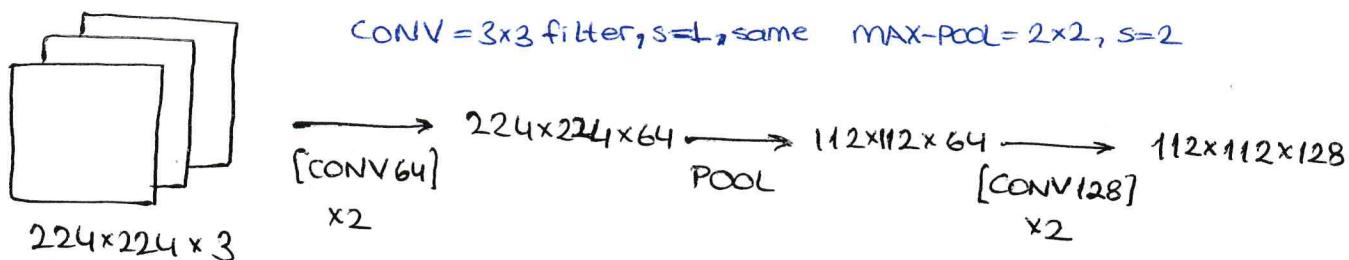
- ReLU

- GPU's were a little bit slower, had a complicated way of training on 2 GPU's (Advanced)

- Local Response normalization (this type of layer isn't used much) (Advanced)



### VGG-16 :



$\sim 138 \text{ M parameters}$

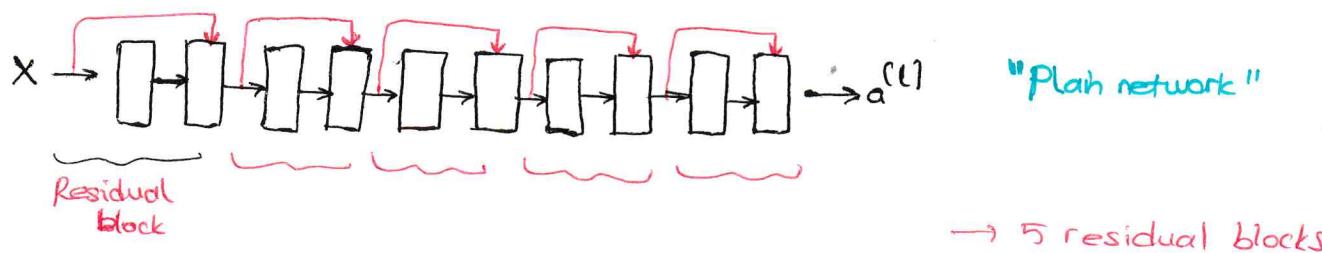
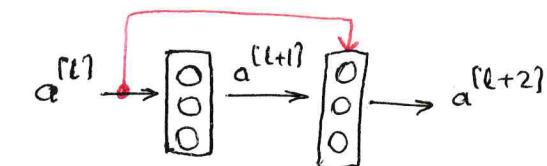
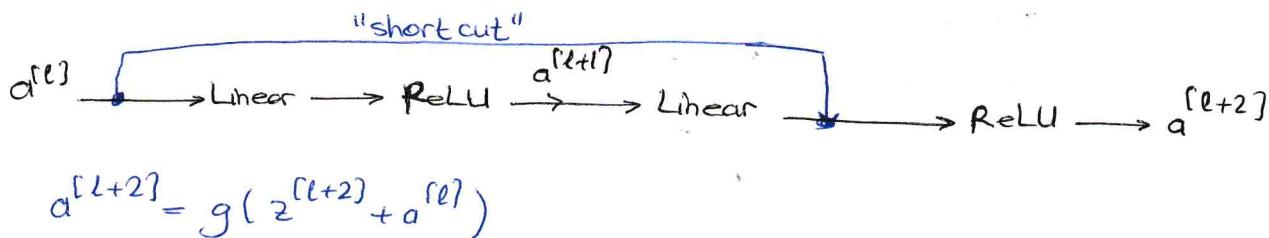
totally 16 layer.

## ResNets

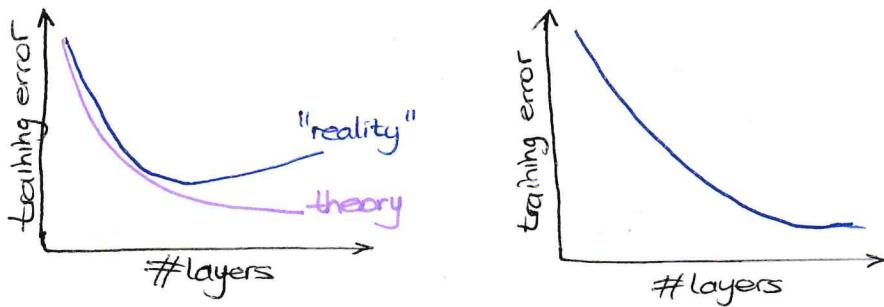
COURSE 4  
WEEK 2

3

Residual Block:



Plain                          ResNet



By taking these activations and allowing it to go much deep in neural network, this really helps with the vanishing and exploding gradient problems and allows to train much deeper networks.

Why Resnets Work,

$$\begin{aligned} a^{l+2} &= g(z^{l+2} + a^l) \\ &= g(W^{l+2}a^{l+1} + b^{l+2} + a^l) \end{aligned}$$

=  $g(a^l)$

If  $W^{l+2} = 0$  and  $b^{l+2} = 0$

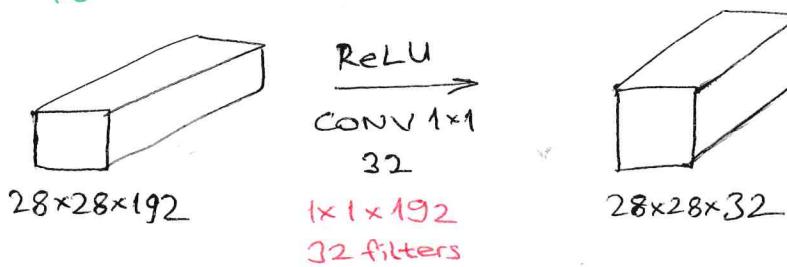
Identity function is easy for residual block to learn! (And it's easy to get  $a^{l+2} = a^l$  because of skip connection)  
We're assuming that  $z^{l+2}$  and  $a^l$  have the same dimension.

## Networks in Networks and 1x1 Convolutions

COURSE 4  
WEEK 2

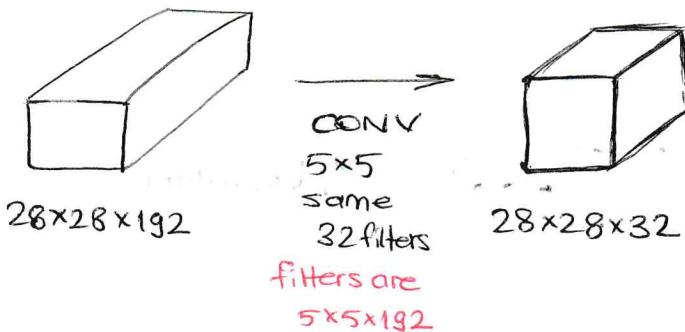
E4

It is sometimes called Network in Network.



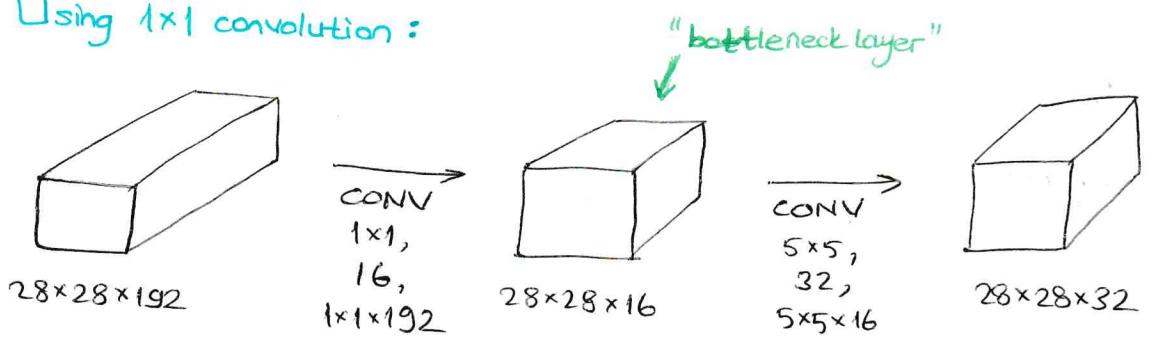
## Inception Network Motivation

Computational cost:



$$28 \times 28 \times 32 \times 5 \times 5 \times 192 \approx 120 \text{ M}$$

Using 1x1 convolution:

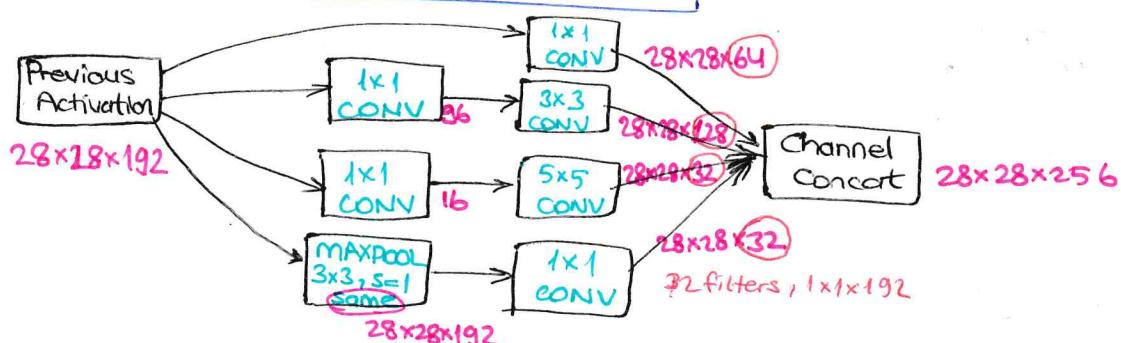


It doesn't seem to hurt the performance but saves you a lot of computation.

$$28 \times 28 \times 16 \times 192 \times 1 \times 1 = 2.4 \text{ M} \quad 28 \times 28 \times 32 \times 5 \times 5 \times 16 = 10.0 \text{ M}$$

12.4 M

## Inception Network



⇒ Inception Block

There can be some side branches.  
FC and softmax  
This appears to have regularizing effect on Inception network and help prevent this network from overfitting.

Transfer Learning

→ If you download ways that someone else has already trained on the network architecture and use that as pre-training and transfer that you might be interested in.

Datasets

ImageNet

MS COCO

Pascal ~~Type~~ of Datasets

# Download some open source implementation of a neural network and download not just the code but also the weights

If you have less data, use pretrained network and with those weights just train your softmax layer.

But, if you have more data, the number of layers you freeze could be smaller and then the number of layers you train on top could be greater.

Data Augmentation

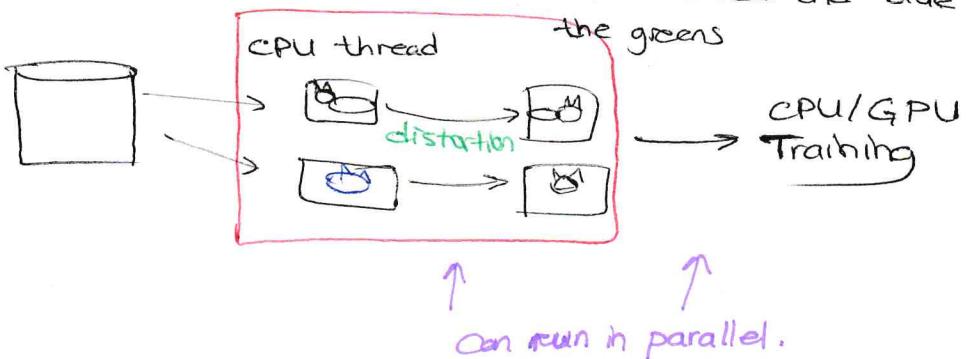
Common Augmentation method :

- \* - Mirroring
- \* - Random cropping
- Rotation
- Shearing
- Local warping
- \* - Color shifting

Advanced !

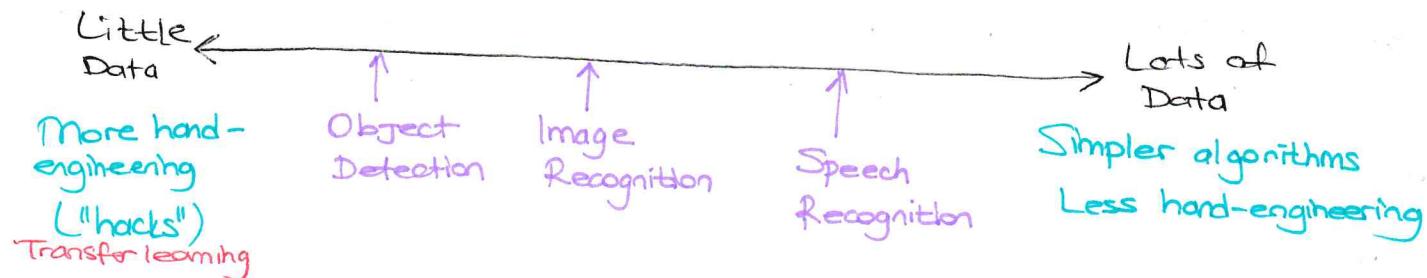
"PCA color augmentation" (in Alexnet paper)

for example, if you have image mainly purple, if it mainly has red and blue tints and very little green, then PCA color augmentation, will add and subtract a lot to red and blue, where it balances relatively to the greens





## Data vs. hard engineering



Two sources of knowledge

- Labelled data
- Hand engineered ~~features~~ / network architecture / other components

When you don't have enough data engineering is very difficult, very skillful task that requires a lot of insight

Tips for doing well on benchmarks / winning competitions :

### Ensembling

- Train several networks independently and average their outputs.  
(This will cause maybe 9/1 or 9/2 better performance)
- (Never use in production)  
unless you have huge computational budget

### Multi-crop at test time

→ (a form of applying data augmentation to your test image)

- Run classifier on multiple versions of test images and average results.  
10-crop (If you have the computational budget, you could use in production)

### Use open source code

- Use architectures of networks published in the literature
- Use open source implementations if possible
- Use pretrained model and fine-tune on your dataset

## FACE RECOGNITION

### One Shot Learning

COURSE 4  
WEEK 4

(1)

When you have ~~one~~ few images. (Such as one)

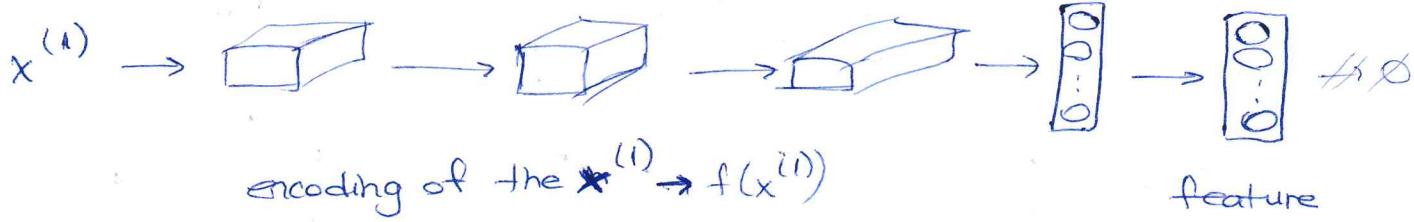
Learning a "similarity" function

$d(\text{img}_1, \text{img}_2)$  = degree of difference between images

If  $d(\text{img}_1, \text{img}_2) \leq \tau$  "same"  
 $> \tau$  "different"

This is how you address the face verification problem

### Siamese Network



encoding of the  $x^{(1)}$   $\rightarrow f(x^{(1)})$

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

**Article**  
**DeepFace** closing the gap to human level performance  
 (Taigmon)

**Summary:** Parameters of NN define an encoding  $f(x^{(i)})$

Learn parameters so that :

If  $x^{(i)}, x^{(j)}$  are the same person,  $\|f(x^{(i)}) - f(x^{(j)})\|^2$  is small

If  $x^{(i)}, x^{(j)}$  are different persons,  $\|f(x^{(i)}) - f(x^{(j)})\|^2$  is large

### Triplet Loss

Anchor	Positive	Negative
A	P	N

$$\underbrace{\|f(A) - f(P)\|^2}_{d(A,P)} + \alpha \leq \underbrace{\|f(A) - f(N)\|^2}_{d(A,N)}$$

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0$$

margin

Loss function :

Given 3 images A, P, N

$$L(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

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

Training set: 10K pictures of 1K 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.

Choose triplets that're hard to train on.

$$d(A, P) + \alpha \leq d(A, N)$$

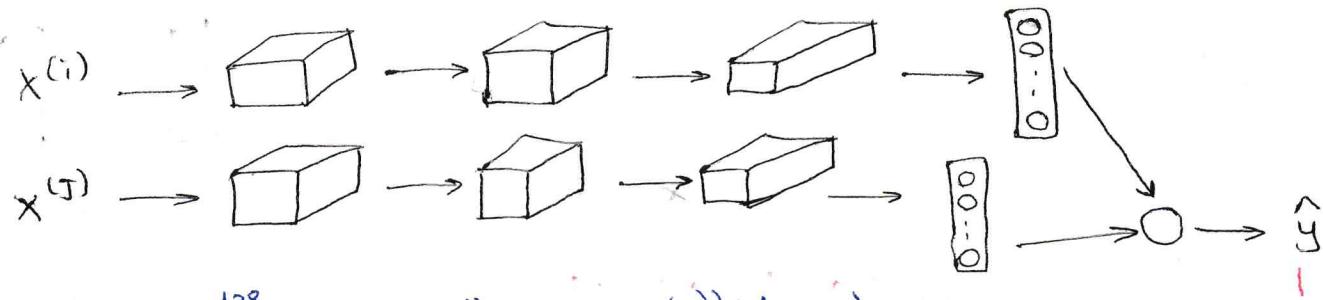
Article X

FaceNet: A unified embedding for face recognition and clustering

Face Net  
Deep Face

} Popular way of naming algorithms in deep learning world.

### Face Verification and Binary Classification



$$\hat{y} = \sigma \left( \sum_{k=1}^{128} w_i |f(x^{(i)})_k - f(x^{(j)})_k| + b \right)$$

$$\frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{f(x^{(i)})_k + f(x^{(j)})_k}$$

$w_i$  and  $b$

from logistic regression unit.

Chi Square

\* Images from database are precomputed ones and when new image come it is compared with precomputed feature vectors of database images.

# Neural Style Transfer

COURSE 4  
WEEK 4

(3)

## What are deep ConvNets Learning ,

Paper by Mathew Zeiler and Rob Fergus

"Visualizing and Understanding Convolutional Networks"

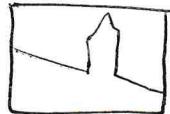
Pick a unit in layer 1. Find the nine image patches that maximize the unit's activation.

Repeat for other units.

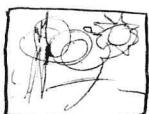
From simple things to complicated things in layers.

Such as edges in Layer 1; to textures in Layer 2, up to detecting very complex objects in the deeper layers .

## Cost Function ,



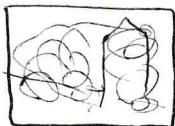
Content C



Style S

$$J(G) = \alpha \cdot J_{\text{content}}(C, G)$$

$$+ \beta \cdot J_{\text{style}}(S, G)$$



Paper by Gatys et al , 2015 . A neural algorithm of artistic style. Images on slide generated by Justin Johnson

## Find the Generated Image G

1. Initiate G randomly.

$$G : 100 \times 100 \times 3$$

RGB

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

$$G := G - \frac{\partial}{\partial G} J$$

## Content Cost Function ,

$$J(G) = \alpha \cdot J_{\text{content}}(C, G) + \beta \cdot J_{\text{style}}(S, G)$$

- Say you use hidden layer  $l$  to compute content cost.
- Use pre-trained ConvNet. (E.g. VGG Network)
- Let  $a^{(l)(C)}$  and  $a^{(l)(G)}$  be the ~~convolution~~ activation of layer  $l$  on the images.
- If  $a^{(l)(C)}$  and  $a^{(l)(G)}$  are similar, both images have similar content.

$$J_{\text{content}}(C, G) = \frac{1}{2} \| a^{(l)(C)} - a^{(l)(G)} \|_2^2$$

not important (Just to normalize but can be adjusted by  $\alpha$ )

→ element-wise sum of squared differences between two activations

## Style Cost Function

COURSE 4  
WEEK 4

(4)

### Style Matrix

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

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

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

in linear algebra this is also called  
"(Gram matrix)", but in this  
video, it is used "style matrix"

$$J_{\text{style}}^{[l]}(S, G) = \frac{1}{(\dots)} \| G^{[l](S)} - G^{[l](G)} \|_F^2$$

$$= \frac{1}{2n_h^{[l]} n_w^{[l]} n_c^{[l]}} \sum_k \sum_{k'} (G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)})^2$$

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

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

$$J(G) = \alpha \cdot J_{content}(C, G) + \beta \cdot J_{style}(S, G)$$

We would like the "generated" image  $G$  to have similar content as the input image  $C$ .

Suppose you have chosen some layer's activations to represent the content of an image.

In practice, you'll get the most visually pleasing results if you choose a layer in the middle of network.

$$a^{rl(C)} \xrightarrow{\text{simply}} a^{(c)} \quad \text{shape: } n_H \times n_W \times n_C$$

$$J_{content}(C, G) = \frac{1}{4 \times n_H \times n_W \times n_C} \cdot \sum_{\text{all entries}} (a^{(c)} - a^{(G)})^2$$

? Each color represents the activations of a filter ~~as~~ it was convolved around an image.

Style Matrix:

Gram matrix  $\rightarrow (v_1, \dots, v_n)$

$$G_{ij} = v_i^T v_j = np.dot(v_i, v_j)$$

$\hookrightarrow$  compares how similar  $v_i$  is to  $v_j$

$$\begin{aligned} \text{default } &\rightarrow \cancel{(n \times n_w \times n_h)} \\ \text{image } &\rightarrow (n_C, n_H \times n_W, n_C) \rightarrow G = (n_C, n_H \times n_W) \cdot (n_H \times n_W, n_C) \\ &\text{result} = (n_C, n_C) \end{aligned}$$

$G_{ij}$  measures how similar the activations of filter  $i$  are to the activations of filter  $j$ .

One important part of the gram matrix is that the diagonal elements such as  $G_{ii}$  also measures how active filter  $i$  is. For example, suppose filter  $i$  is detecting vertical textures in the image. Then  $G_{ii}$  measures how common vertical textures are in the image as a whole; if  $G_{ii}$  is large, this means that the image has a lot of vertical texture.

STYLE WEIGHTS

$$J_{\text{style}}(s, g) = \sum_{\ell} \lambda^{(\ell)} J_{\text{style}}^{(\ell)}(s, g)$$