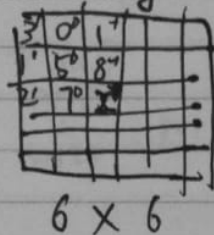


Coursera Deep learning course 4 Note

CNN

Week 1

Vertical edge detection



6 x 6

"convolution"

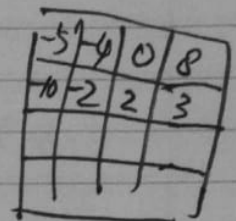


3 x 3

*

filter

=



4 x 4

$$3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0$$

$$+ 5 \times 0 + 7 \times 0 + 1 \times -1 + 8 \times -1 + 2 \times -1 = -5$$

10 10 10 0 0 0
10 10 10 0 0 0
10 10 10 0 0 0
10 10 10 0 0 0
10 10 10 0 0 0
10 10 10 0 0 0

6 x 6

* $\begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix}$

3 x 3

=

$\begin{matrix} 0 & 3 & 0 & 3 & 0 & 0 \\ 0 & 3 & 0 & 3 & 0 & 0 \\ 0 & 3 & 0 & 3 & 0 & 0 \\ 0 & 3 & 0 & 3 & 0 & 0 \end{matrix}$

4 x 4

↓ vertical edges

$\begin{matrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{matrix}$

sobel filter

$\begin{matrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{matrix}$

schant filter

$\left. \begin{matrix} w_1, w_2, w_3 \\ w_4, w_5, w_6 \\ w_7, w_8, w_9 \end{matrix} \right\} \text{ can be trained}$

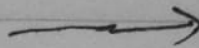
Padding

6x6
n x n

3x3
f x f

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

$$6 - 3 + 1 = 4$$



4x4

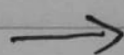
Padding

6x6 → 8x8
n x n n+2p
 x
 n+2p

3x3
f x f

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

$$8 - 3 + 1 = 6$$



6x6

"Valid" convolution $n \times n * f \times f \rightarrow n - f + 1 \times n - f + 1$

usually odd



"Same" convolution $(n + 2p - f + 1) \times (n + 2p - f + 1)$

Strided convolution

7x7

* 3x3

= 3x3

↓ floor

stride = 2

f x f

n x n

*

p=0 s=2

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

3x3

Convolution on RGB images

$$\begin{matrix} 6 \times 6 \times 3 \\ \text{height width channels} \end{matrix} * \begin{matrix} 3 \times 3 \times 3 \\ \text{height width channels} \end{matrix} = \begin{matrix} 4 \times 4 \\ \text{height width} \end{matrix}$$

R, G, B

$$\begin{matrix} 27 \text{ numbers} \\ (9+9+9) \end{matrix} * \begin{matrix} 27 \text{ numbers} \\ R \quad G \quad B \end{matrix} = 1 \text{ number}$$

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

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

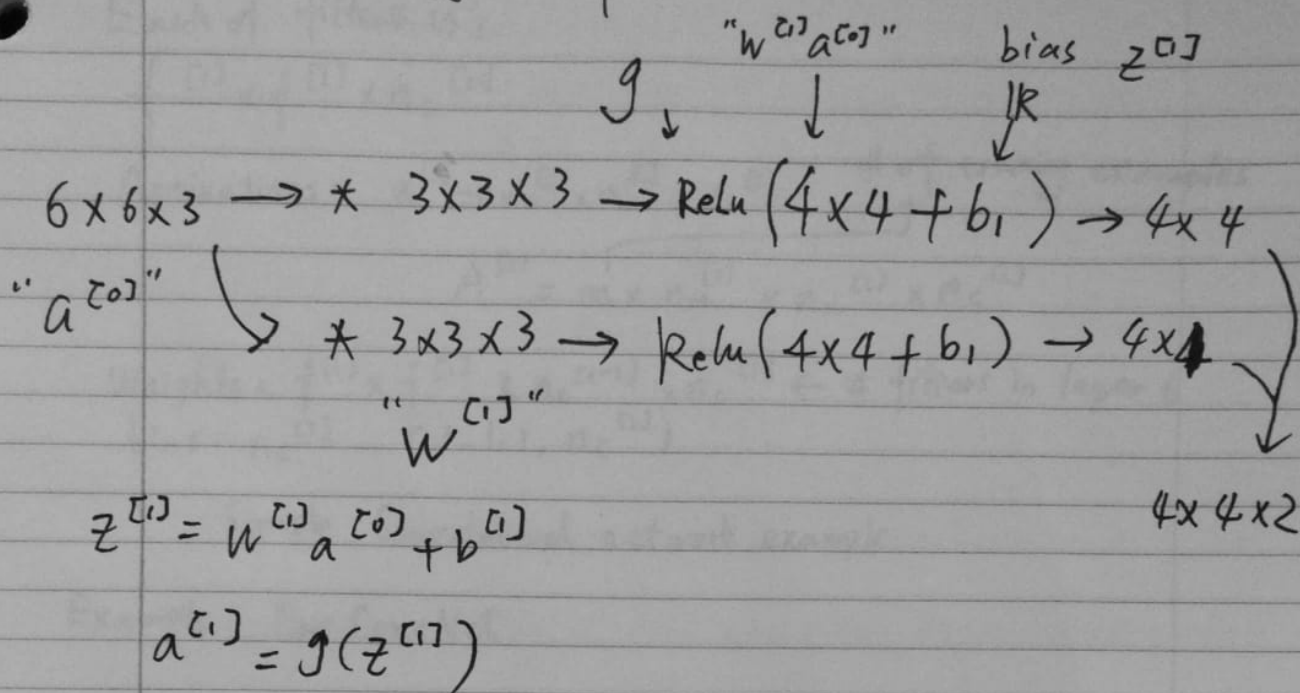
$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$

multiple filter

$$\begin{matrix} 6 \times 6 \times 3 \end{matrix} \begin{matrix} \nearrow * \\ \searrow * \end{matrix} \begin{matrix} \text{vertical edge} \\ 3 \times 3 \times 3 = 4 \times 4 \\ \text{horizontal edge} \\ 3 \times 3 \times 3 = 4 \times 4 \end{matrix} \rightarrow 4 \times 4 \times 2$$

$$\begin{matrix} n \times n \times n_c \\ \uparrow \\ \text{channel} \\ \text{(or depth)} \end{matrix} * \begin{matrix} f \times f \times f \\ \# n_c' \text{ filters} \end{matrix} \rightarrow \begin{matrix} n-f+1 \times n-f+1 \times n_c' \\ \uparrow \\ \# n_c' \text{ filters} \end{matrix}$$

One layer of convolution network



e.g.: 10 filters that are $3 \times 3 \times 3$ in one layer
How many parameters does that layer have?

$3 \times 3 \times 3$
27 parameters
+ bias
= 28 parameters $\times 10 = 280$ parameters

Summary of notation

$f^{[l]}$ = filter size
 l = layer number
 $s^{[l]}$ = stride
 $n_f^{[l]}$ = num of filters

Input: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$
Output: $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$n^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2P^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$

Each of filters is:

$$f^{[l]} \times f^{[l]} \times n_c^{[l]}$$

Activations: $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$ # of training examples

$$A^{[l]} = m \times n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$$

Weights: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]} \leftarrow \# \text{ filters in layer } l$
 bias: $n_c^{[l]} - (1, 1, 1, n_c^{[l]})$

Simple Convolutional network example

Example ConvNet

$39 \times 39 \times 3$	$\xrightarrow{f^{[1]}=3}$	$37 \times 37 \times 10$	$\xrightarrow{f^{[2]}=5}$	$17 \times 17 \times 20$
$n_H^{[0]} = n_W^{[0]} = 39$	$s^{[1]}=1$	$n_H^{[1]} = n_W^{[1]} = 37$	$s^{[2]}=2$	$n_H^{[2]} = n_W^{[2]} = 17$
$n_c^{[0]} = 3$	$p^{[1]}=0$	$n_c^{[1]} = 10$	$p^{[2]}=0$	$n_c^{[2]} = 20$
	$n_c^{[1]} = 10$		$n_c^{[2]} = 20$	

$\rightarrow *$	$7 \times 7 \times 40$	$\xrightarrow{\text{softmax}}$	\hat{y}
$f^{[3]}=5$	$n_H^{[3]} = n_W^{[3]} = 7$		
$s^{[3]}=2$	$n_c^{[3]} = 40$	$1960 \rightarrow 1$	
$n_c^{[3]} = 40$			

$n_H, n_W \downarrow$

$n_c \uparrow$

Types of layer

Convolution (CONV)

Pooling (POOL)

Fully Connected (FC)

Pooling layers

Max pooling: take max of numbers in a region

$$4 \times 4 \rightarrow 2 \times 2$$

$f=2$
 $s=2 \rightarrow$ parameters do not change in gradient decent

$$5 \times 5 \times 2 \xrightarrow[f=3]{s=1} 3 \times 3 \times 2$$

Average pooling: take average of numbers in a region

note: max pooling is used more often than average pooling

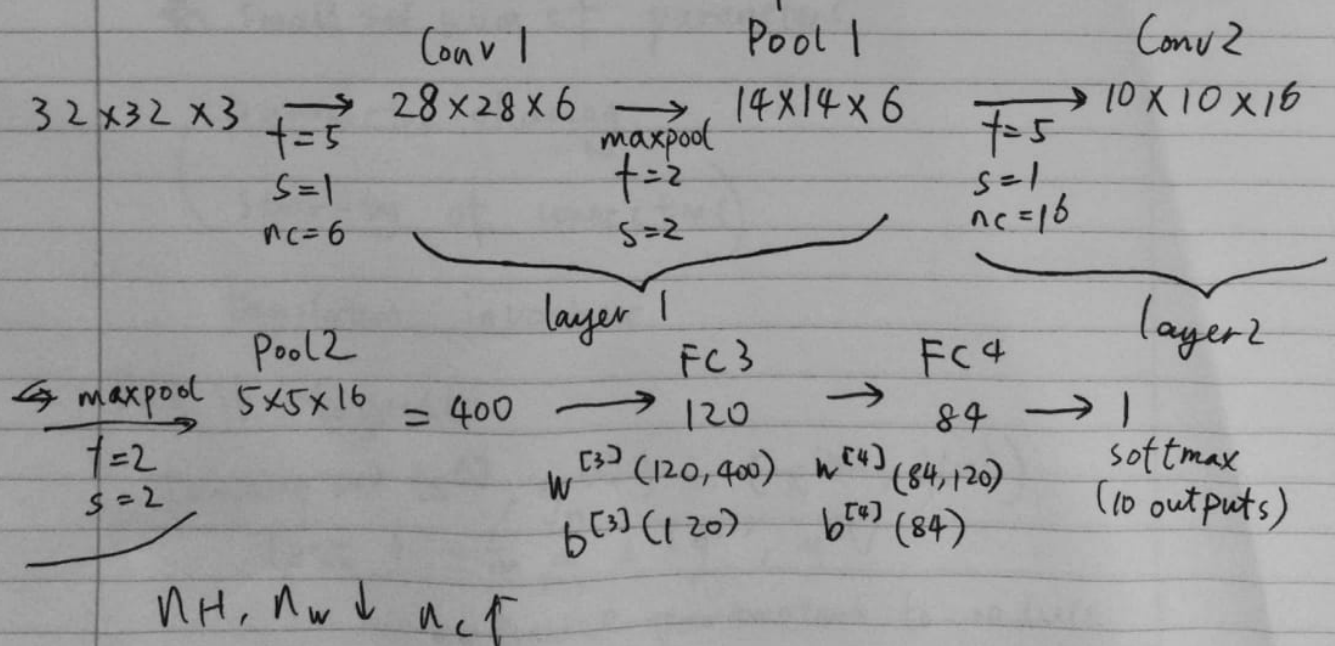
Hyperparameters

f , s , max or average pooling

No parameters to learn

$$\begin{aligned} & n_H \times n_W \times n_C \\ & \downarrow \\ & \left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor \\ & \times n_C \end{aligned}$$

Neural network example (Le Net - 5)



	Shape	Activation shape	Activation Size	# parameters
Input		$(32, 32, 3)$	$32 \times 32 \times 3$	
		\downarrow	3072	$a^{(0)}$
CONV1 ($f=5, s=1$)		$\frac{32+0-5}{1}+1$ \downarrow $(28, 28, 8)$	$28 \times 28 \times 8$	$5 \times 5 \times 8 + 8$
		\downarrow	6272	\downarrow 208
POOL1		$28/2$ \downarrow $(14, 14, 8)$	$14 \times 14 \times 8$	no parameters to learn
		\downarrow	1568	\downarrow 0
CONV2 ($f=5, s=1$)		$\frac{14+0-5}{1}+1$ \downarrow $(10, 10, 16)$	$10 \times 10 \times 16$	$5 \times 5 \times 16 + 16$
		\downarrow	1600	\downarrow 416
POOL2		$10/2$ \downarrow $5 \times 5 \times 16$	$5 \times 5 \times 16$	no parameters to learn
		\downarrow	400	\downarrow 0
FC3		$(120, 1)$	120	$400 \times 120 + 1$
				\downarrow 48001
FC4		$(84, 1)$	84	$120 \times 84 + 1$
				\downarrow 10081
Softmax		$(10, 1)$	10	$84 \times 10 + 1$
				\downarrow 841

Why convolutions

* Small set num of parameters

(parameter sharing
sparsity of connections)

translation invariance

Put it together

Training set $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$

$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

minimize parameters to reduce

Classic networks:

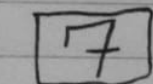
- LeNet-5
- AlexNet
- VGG

ResNet (152)

Inception

Case studies
Classic networks

LeNet-5



$32 \times 32 \times 1$

$\xrightarrow[5 \times 5]{s=1}$

$28 \times 28 \times 6$

$\xrightarrow[5 \times 5]{f=2, s=2}$

$14 \times 14 \times 6$

$\xrightarrow[5 \times 5]{s=1}$

$10 \times 10 \times 6$

$\xrightarrow[5 \times 5]{f=2, s=2}$

$5 \times 5 \times 16$
 $\underline{400}$

\xrightarrow{FC}

120

\xrightarrow{FC}

84

$\rightarrow \hat{y}$
Softmax (10)

60k parameters

$n_H, n_W \downarrow \quad n_C \uparrow$

conv pool conv pool fc fc output

AlexNet

$$227 \times 227 \times 3 \xrightarrow[11 \times 11]{s=4} 55 \times 55 \times 96 \xrightarrow[3 \times 3]{\text{max-pool}, s=2} 27 \times 27 \times 96 \xrightarrow[5 \times 5]{\text{same}} 27 \times 27 \times 256$$

$$\xrightarrow[3 \times 3]{\text{max-pool}, s=2} 13 \times 13 \times 256 \xrightarrow[3 \times 3]{\text{same}} 13 \times 13 \times 384 \xrightarrow[3 \times 3]{\text{same}} 13 \times 13 \times 384 \xrightarrow[3 \times 3]{\text{same}} 13 \times 13 \times 256$$

$$\xrightarrow[3 \times 3]{\text{max-pool}, s=2} 6 \times 6 \times 256 = 9216 \xrightarrow{\text{FC}} 4096 \xrightarrow{\text{FC}} 4096 \xrightarrow{\text{softmax}} 1000$$

- Similar to LeNet, but much bigger: 60 m parameters

VGG-16

CONV = 3×3 filters, $s=1$, same MAX-POOL = 2×2 , $s=2$

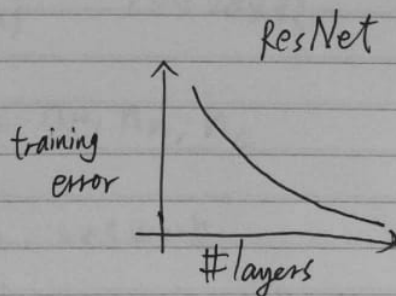
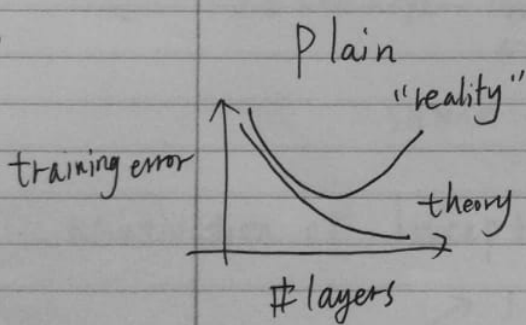
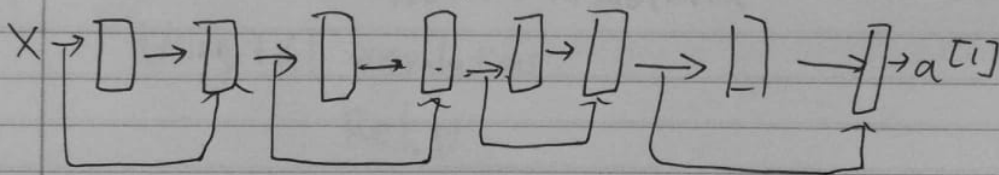
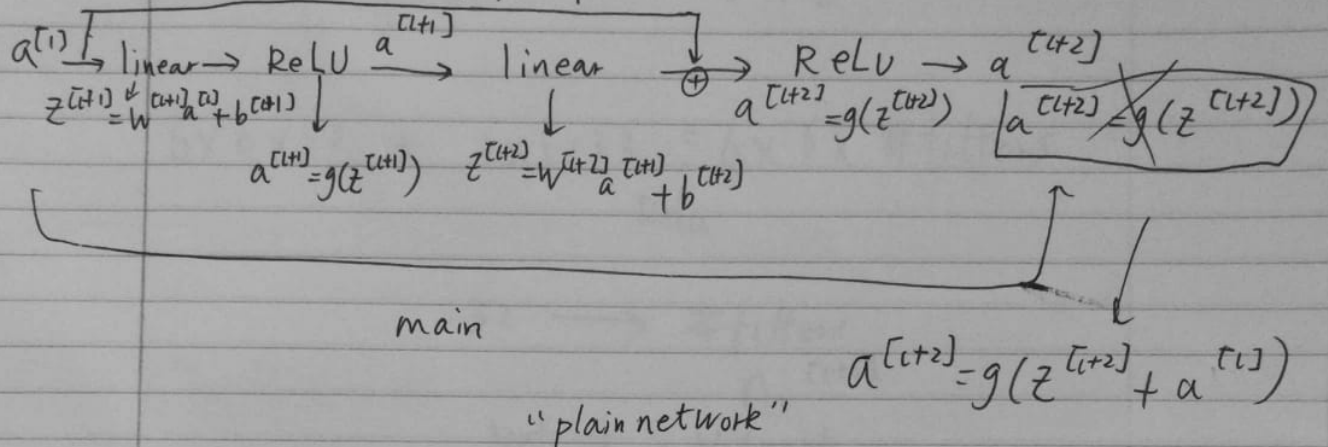
$$224 \times 224 \times 3 \xrightarrow[\text{CONV } 64 \times 2]{} 224 \times 224 \times 64 \xrightarrow{\text{POOL}} 112 \times 112 \times 64 \xrightarrow[\text{CONV } 128 \times 2]{} 112 \times 112 \times 128$$

$$\xrightarrow{\text{POOL}} 56 \times 56 \times 128 \xrightarrow[\text{CONV } 256 \times 3]{} 56 \times 56 \times 256 \xrightarrow{\text{POOL}} 28 \times 28 \times 256 \xrightarrow[\text{CONV } 512 \times 3]{} 28 \times 28 \times 512$$

$$\xrightarrow{\text{POOL}} 14 \times 14 \times 512 \xrightarrow[\text{CONV } 512 \times 3]{} 14 \times 14 \times 512 \xrightarrow{\text{POOL}} 7 \times 7 \times 512 \xrightarrow{4096} \text{FC} \xrightarrow{4096} \text{FC} \xrightarrow{1000} \text{softmax}$$

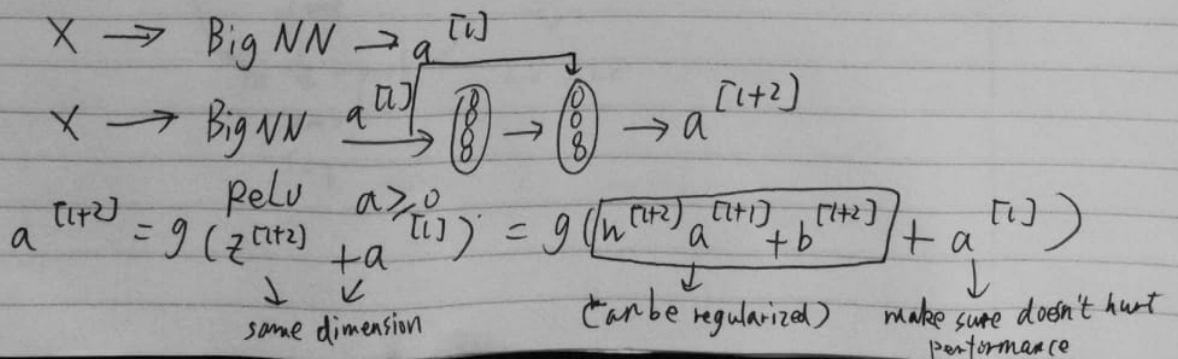
weighted layers
16 weights
(38 m parameters)

Residual block (help train very deep network)
 "short cut" / skip connection



Why Resnet works?

can be also $g(z^{[l+2]} + b^{[l+2]} + w_s a^{[1]})$
 help adjust dimension



1×1 convolution

$$6 \times 6 * \boxed{\begin{matrix} \diagup & \diagdown \\ \diagdown & \diagup \end{matrix}} 1 \times 1 = 6 \times 6$$

$$6 \times 6 \times 32 * 1 \times 1 \times 32 = 6 \times 6 \times \# \text{filters}$$

Relu

$$32 \longrightarrow \# \text{filters}$$

$n_c^{(l+1)}$

Network in Network

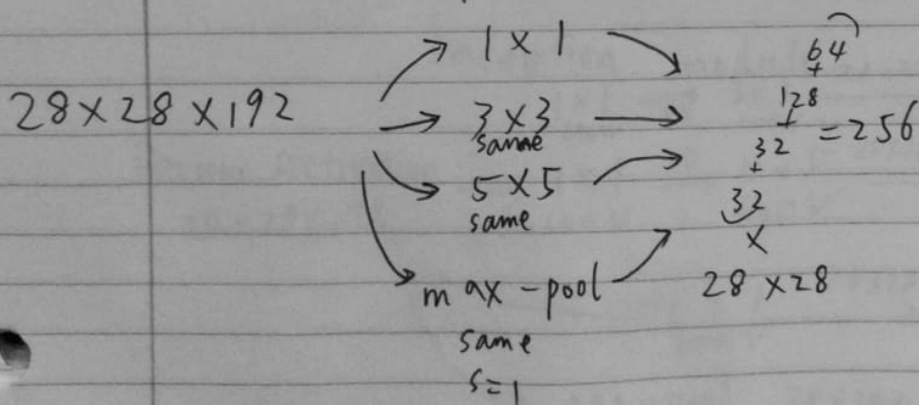
Using 1×1 convolutions

ReLU

$$28 \times 28 \times 192 \xrightarrow[\substack{\text{conv } 1 \times 1 \\ 32}]{\text{ReLU}} 28 \times 28 \times 32$$

shrink to n_H, n_W, n_C

Motivation of Inception network



The problem of computational cost

$$28 \times 28 \times 192 \xrightarrow[\substack{\text{conv} \\ 5 \times 5, \\ \text{same} \\ 32}]{ } 28 \times 28 \times 32$$

32 filter

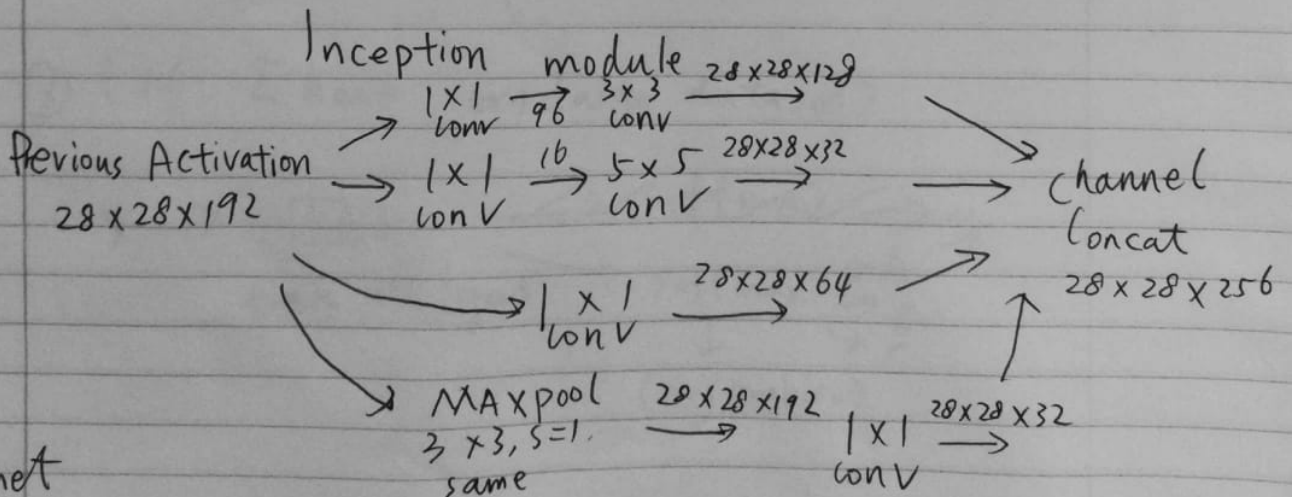
Computation: $28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120 \text{ M}$

Using 1×1 convolution "bottleneck layer"

$$28 \times 28 \times 192 \xrightarrow[\substack{\text{conv} \\ 1 \times 1, \\ 16, \\ 1 \times 1 \times 192}]{ } 28 \times 28 \times 16 \xrightarrow[\substack{\text{conv} \\ 5 \times 5, \\ 32, \\ 5 \times 5 \times 16}]{ } 28 \times 28 \times 32$$

Computation: $28 \times 28 \times 16 \times 192 = 2.4 \text{ m}$

+ $28 \times 28 \times 32 \times 5 \times 5 \times 16 = 10.0 \text{ m}$ = 12.4 m



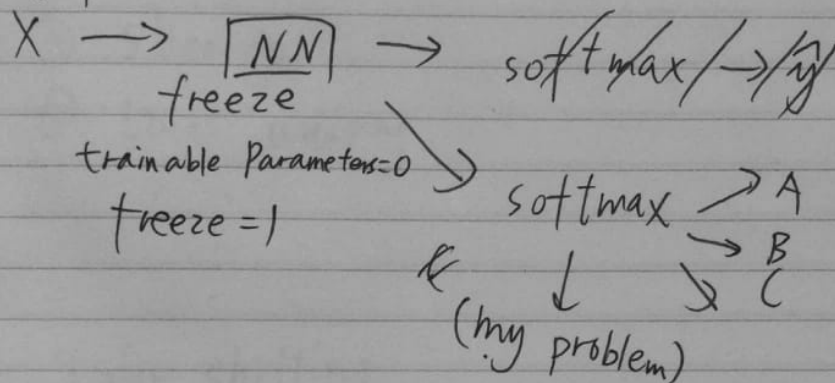
GoogLeNet

advice for using convNets

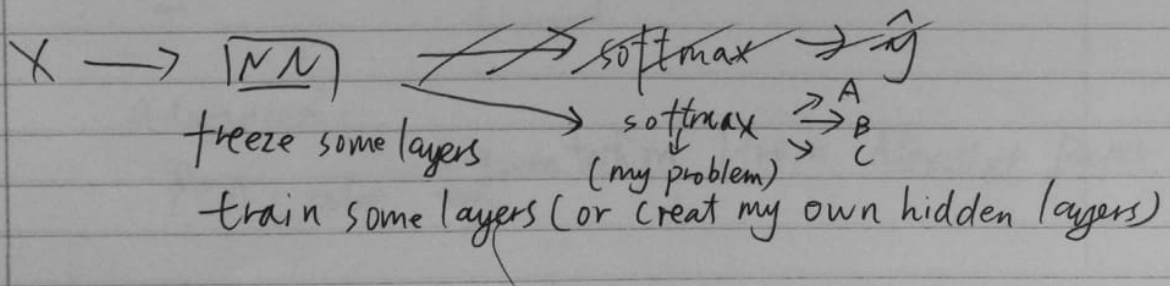
1. Using open-source implementation

2. transfer learning

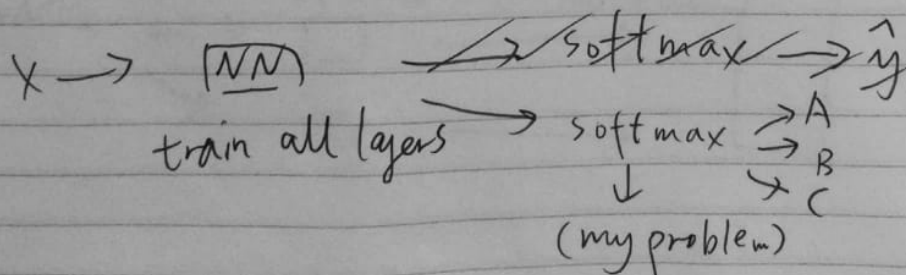
① (if I have small dataset)



② (if I have many pictures)



③ (if I have very large dataset)



3. Data augmentation

① mirroring

② Random Cropping

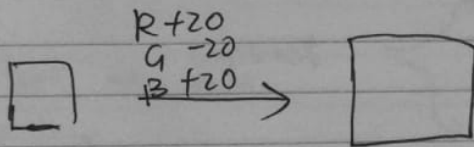
③ Rotation

④ shearing

⑤ Local warping

⋮

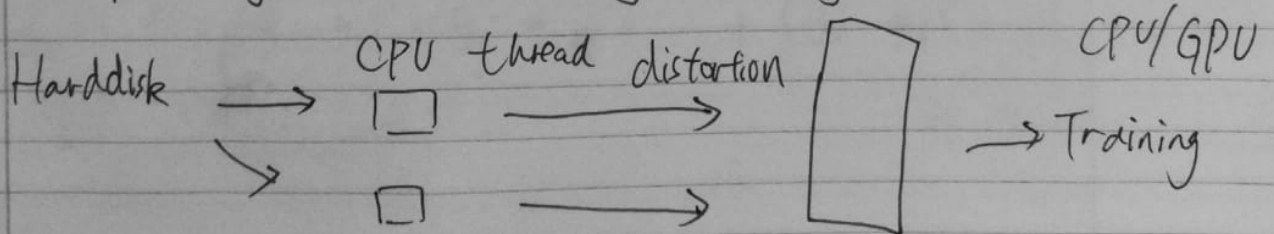
Color shifting



Advanced:

PCA color augmentation (from AlexNet Paper)

Implementing distortions during training

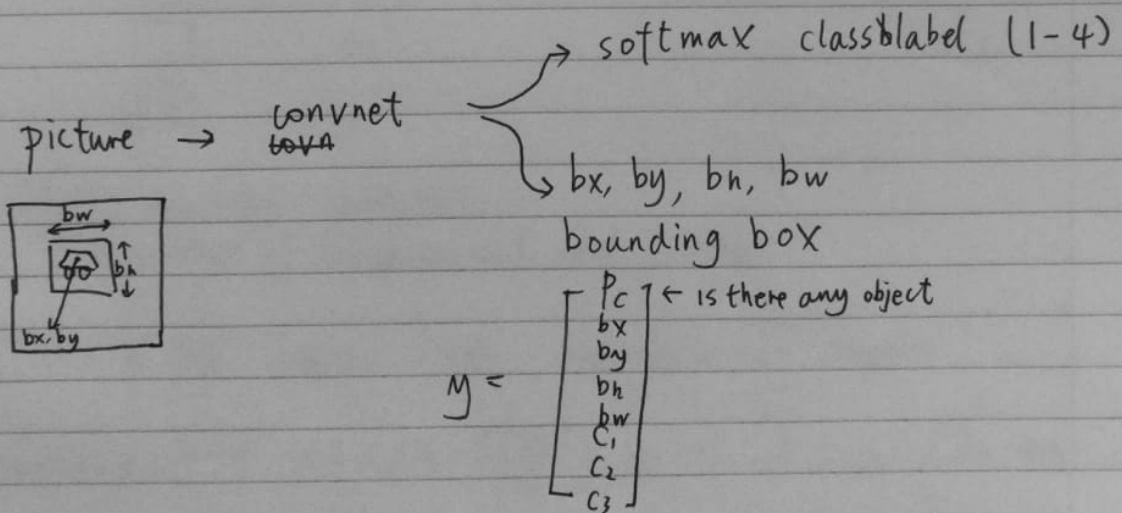


Tips for doing well on benchmarks
 Ensembling (not used in production)
 (Average different models' output)
~~Multi~~

Multi-crop at test time

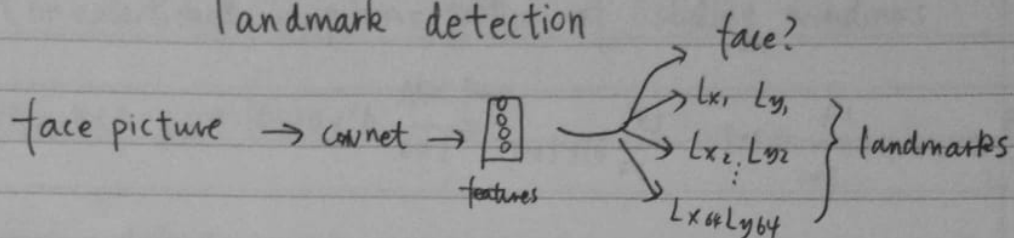
Classification with localization

1. pedestrian
2. car
3. motorcycle
4. background



$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2 & \text{if } y_1 = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0 \end{cases}$$

landmark detection



e.g



• black points are landmarks

Object detection

Car detection example

Training set



1



1



0



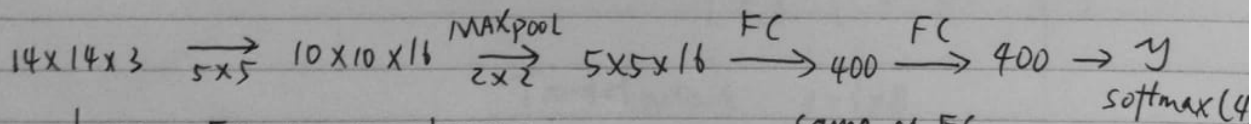
1

\rightarrow convnet

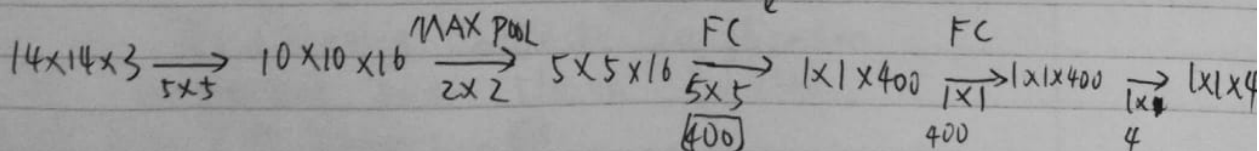
Sliding windows detection

Disadvantage: computational cost, slow

Turning FC layer into convolutional layers



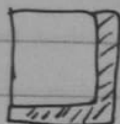
change FC to conv layer:



Convolutional implementation of sliding windows

$$14 \times 14 \times 3 \xrightarrow{5 \times 5} 10 \times 10 \times 16 \xrightarrow[2 \times 2]{\text{MAX Pool}} 5 \times 5 \times 16 \xrightarrow[5 \times 5]{\text{FC}} 1 \times 1 \times 400 \xrightarrow[1 \times 1]{\text{FC}} 1 \times 1 \times 400 \xrightarrow[1 \times 1]{\text{FC}} 1 \times 1 \times 4$$

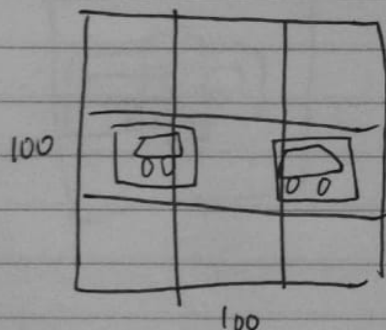
add stripe as follows



$$16 \times 16 \times 3 \xrightarrow{5 \times 5} 12 \times 12 \times 16 \xrightarrow[2 \times 2]{\text{MAX Pool}} 6 \times 6 \times 16 \xrightarrow[5 \times 5]{\text{FC}} 2 \times 2 \times 400 \xrightarrow[1 \times 1]{\text{FC}} 2 \times 2 \times 400 \xrightarrow[1 \times 1]{\text{FC}} 2 \times 2 \times 4$$

disadvantage: bounding box not too accurate

Output accurate bounding box
YOLO algorithm



Labels for training
For each grid cell

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

Target output $3 \times 3 \times 8$

Specify the bounding boxes

Evaluating object localization

Intersection over Union (IOU)

= size of intersection

Size of

"correct" if $\text{IOU} \geq 0.5$



Non-max suppression example

look at the probability of detections, and output max one, suppress lower probability ones

Discard all boxes with $p_c \leq 0.6$

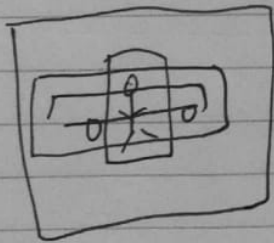
While there are many remaining boxes

- Pick the box with largest p_c

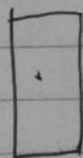
Output that as a prediction

- Discard any remaining box with $IOU \geq 0.5$ with the box output in the previous step

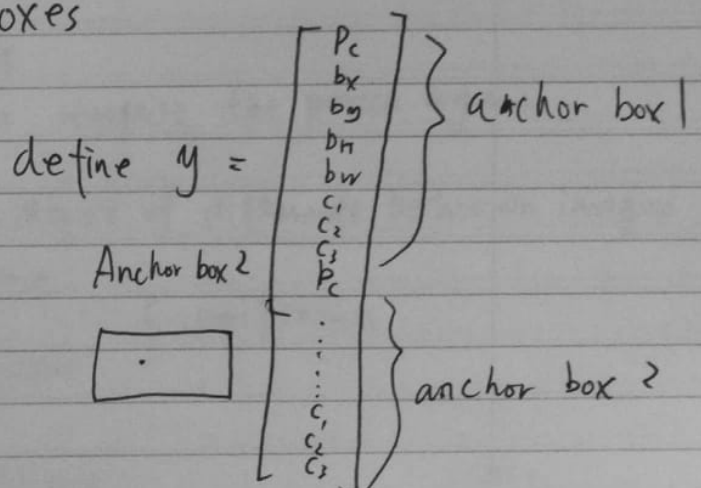
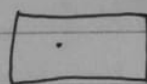
Anchor boxes



Anchor box 1



Anchor box 2



object is assigned to grid cell that contains object's mid point and anchor box which has highest IOU

Region proposal
R-CNN

Propose regions

of sliding window

Fast R-CNN

Use convolution implementation to ~~propose~~ ^{✓ propose} regions

Faster R-CNN

Use convolution network to propose regions

Face verification vs. face recognition

Verification

- Input image, name/ID
- Output whether the input image is that of the 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")

One-shot learning

Learn from one example to recognize the person again

$d(\text{img1}, \text{img2}) = \text{degree of difference between images}$
if $\begin{cases} d \leq \tau & \text{"same"} \\ > \tau & \text{"different"} \end{cases} \text{ verification}$

Siamese network

$x^{(1)} \rightarrow \text{CNN} \rightarrow \cancel{\text{softmax}} f(x^{(1)})$ "encoding of $x^{(1)}$ "

$x^{(2)} \rightarrow \text{CNN} \rightarrow f(x^{(2)})$ "encoding of $x^{(2)}$ "

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

Learning Objective

(A) (P)
Anchor Positive
d to be large small

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

(A) (N)
Anchor negative
d to be large

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

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

Loss function

margin \rightarrow make sure
 $f(\text{img}) = 0$
condition is covered

Given 3 images A, P, N

$$\mathcal{L}(A, P, N) = \max(\underbrace{\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha}_{> 0}, 0)$$

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

Training set: 10k pictures of 1k persons

(Note: there may be more than 1 pics for one person)

Choosing the triplets of A, P, N

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

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

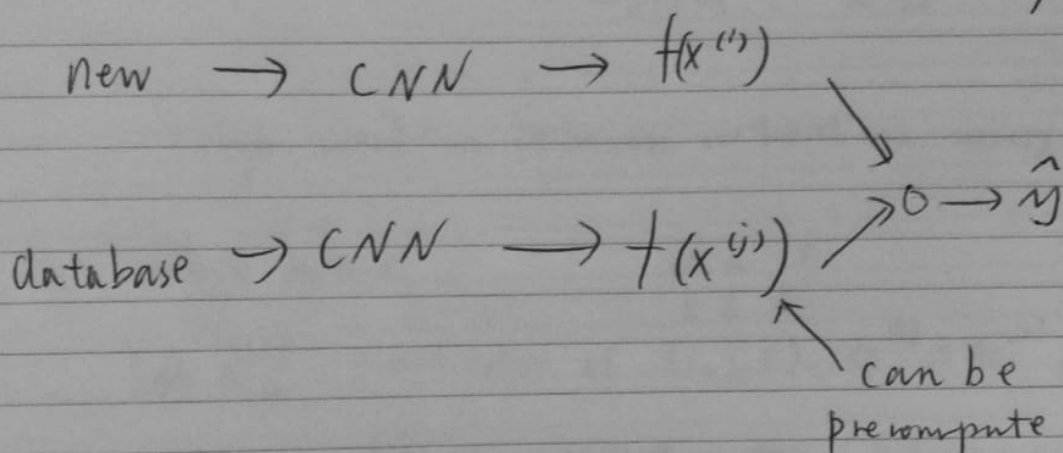
$$d(A, P) \approx d(A, N)$$

Verification (as binary classification)

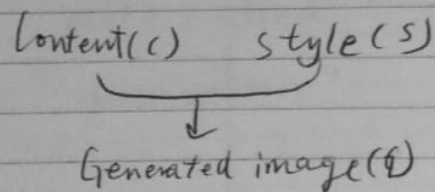
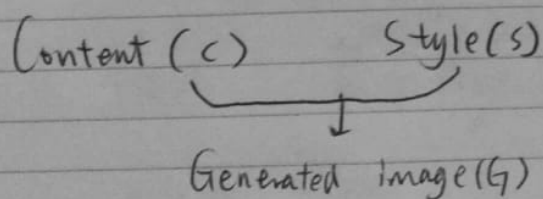
Learning the similarity function

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

$$\text{or} \quad \left(\frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{f(x^{(i)})_k + f(x^{(j)})_k} \right) \times x^2$$



Neural style transfer



Loss function

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

Find the generated image G

1. Initiate G randomly
 $G: 100 \times 100 \times 3$
2. Use gradient descent to minimize $J(G)$

$$G = G - \frac{\alpha}{2G} J(G)$$

Content cost function

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

choose: ~~choose~~ layer L : $a^{[L]}(C)$, $a^{[L]}(G)$ are the activations

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

style: correlation between activations across channels

style matrix

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

$G_{kk'}^{[L]} \leftarrow$ compare correlations between k and k'

\leftarrow style image

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

also called "gram matrix"

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

$$J_{\text{style}}^{[L]}(S, G) = \frac{1}{(2n_H^{[L]}n_W^{[L]}n_C^{[L]})} \|G^{[L]}(S) - G^{[L]}(G)\|_F^2$$

$$= \frac{1}{(2n_H^{[L]}n_W^{[L]}n_C^{[L]})} \sum_k \sum_{k'} (G_{kk'}^{[L]}(S) - G_{kk'}^{[L]}(G))^2$$

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

Convolutions in 2D and 1D

2D

$$14 \times 14 \times 3 * 5 \times 5 \times 3 \rightarrow 10 \times 10 \times 16$$

1D

$$14 \times 1 * 5 \times 1 \rightarrow 10 \times 16$$

3D

$$\begin{matrix} H & W & D & \downarrow & n_c \\ 14 \times 14 \times 14 \times 1 * 5 \times 5 \times 5 \times 1 \end{matrix} \rightarrow \begin{matrix} 10 \times 10 \times 10 \times 16 \\ n_{c+1} \end{matrix}$$

16 filter

n_{c+1}