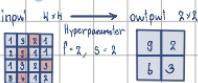
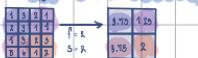


reduce the size of the inputs

11 Pooling layer : Max pooling



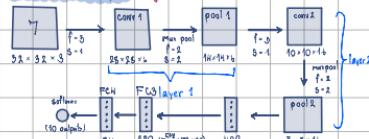
12 Pooling layer : Average pooling



Summary of pooling

Hyperparameters: f : filter size
 s : stride

13. CNN examples e.g. LeNet - 5



14 Why convolutions



train as a vertical edge detector
Parameter sharing → A feature detector

2) Sparsity of connections → In each layer, back output value depends only on a small number of inputs

Putting it together

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



$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m h(x^{(i)}) - y^{(i)} \|^2$$

use gradient descent to optimize reduce J

15. Why look at case studies

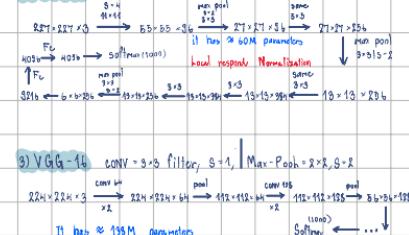
- Convolutional networks → LeNet-5, AlexNet, VGG
- ResNet (152 layer)
- Inception

16. Classic Networks

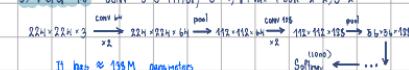
1) LeNet - 5



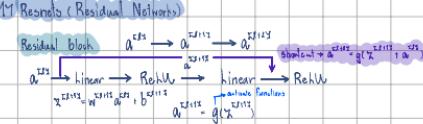
2) Alex net



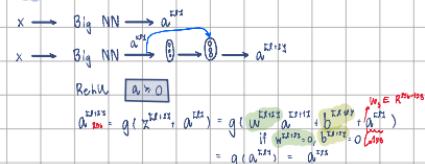
3) VGG - 16



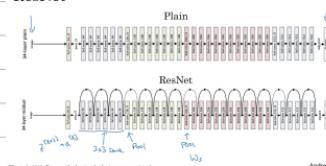
17. Resnets (Residual Networks)



18. Why ResNets work



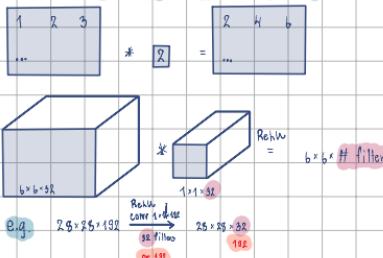
ResNet



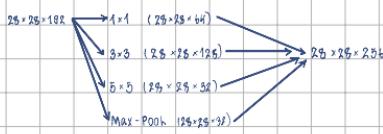
[He et al., 2015]. Deep residual networks for image recognition

Andrew Ng

19 Network in Network and 1×1 convolutions



20 Inception Network Motivation



The problem of computation cost

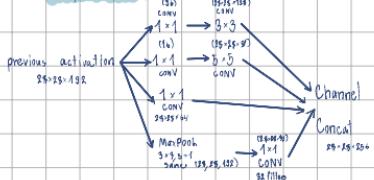
$$\begin{aligned} & 25 \times 25 \times 192 \xrightarrow{\text{conv}} 25 \times 25 \times 32 \\ & (25 \times 25 \times 32) \times (5 \times 5 \times 192) = 120 \text{ million} \end{aligned}$$

just using 1×1 convolution

$$\begin{aligned} & 25 \times 25 \times 192 \xrightarrow{\text{conv}} 25 \times 25 \times 32 \\ & [25 \times 25 \times 192 \cdot 1 \times 1 \times 192] = 104 \text{M} = 12.4 \text{M} \end{aligned}$$

21 Inception Network (put a lot of these modules together)

Inception module



Inception network \Rightarrow a lot of Inception module put together

22 Transfer learning

e.g. 3 classes \rightarrow Trigger \rightarrow Misfit \rightarrow Neither

download implementation NN, change softmax to match your data
just train softmax for outputs or freeze some layers and train
 $x \rightarrow \dots \rightarrow \text{softmax} \rightarrow y$

23 Data Augmentation \rightarrow improve performance computer vision

Common augmentation method: 1) Mirroring $\square \rightarrow \square$

2) Random Cropping

3) Rotation

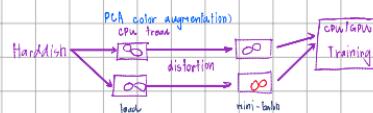
4) Shearing

5) local warping...

:

6) Color shifting \rightarrow increase or decrease some color channel

Advanced: PCA, Min class avg, AlexNet paper



24 State of computer vision

Data vs hand engineering

little data \rightarrow transfer learning \rightarrow lots of work

Two source of knowledge

• labeled data (x, y) \rightarrow needs

• Human engineer feature, network architecture, other comp. tips - Assembling \rightarrow train several networks independently and average their outputs

- Multi-crop at test time \rightarrow run classifier on multiple versions of test images result

use open source code implementation

use pretrained model and fine-tune on your dataset

Object Detection

2.5 Object localization

Image classification \rightarrow input picture, then output result [image]

Classification with localization \rightarrow putting bounding box around result [image]

Detection \rightarrow multiple objects [B]

Classification with localization



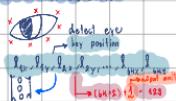
Define the target label y

need to output b_x, b_y, b_w, b_h , class label

$$y = \begin{cases} p_0 & \text{if there is an object (probability)} \\ 0 & \text{else} \end{cases}$$

$$\begin{cases} b_x, b_y, b_w, b_h & \text{if } y=1 \\ (x - \bar{x})^2 + (y - \bar{y})^2 & \text{if } y=0 \end{cases}$$

2.6 Handmark Detection



2.7 Object Detection

example training set = (X, Y)

image \rightarrow ConvNet \rightarrow y

Sliding windows detection



change size of window

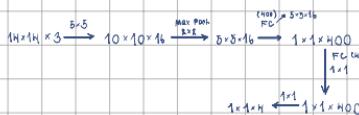
from small to big (every iteration)

disadvantage \rightarrow computation cost (use more resources)

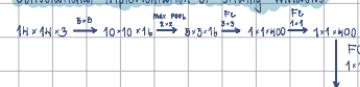
\rightarrow better performance than sliding windows detection

2.8 Convolutional implementation of sliding windows

Turning FC layer into convolutional layer



Convolutional implementation of sliding windows



2.9 Intersection Over Union (How to tell obj algorithm is working well)

Evaluating object localization calculate intersection over union



= size of \cap

size of \cup

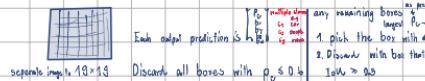
"correct" if IoU ≥ 0.5

3.0 Non-Max Suppression (make sure algorithm detect only one)

e.g. p_i (probability of detection)

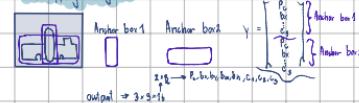


p_i vs p_j



3.1 Anchor Boxes

Overlapping objects



3.2 Yolo Algorithm



Masking prediction P \rightarrow number of classes

$O_{i,j,k} \rightarrow$ check classes

Outputting the non-max suppressed output

- for each grid cell, get 2 predicted bounding boxes
- get rid of low probability predictions
- for each class use non-max suppression to generate final prediction

33. Region Proposals (R-CNN)

- pick just few regions that make sense
- to run (instead running sliding windows)



Segmentation Algorithm

For R-CNN → use convolution implementation.

Faster R-CNN
use convolution
+ network of sliding windows to classify all the proposed regions

34 What is face recognition?

Face Verification	Face recognition
- Input image, name/ID	new answer > 21
- Output whether the input images in that of the claimed person	- Has a database of K persons
	- Get an input image
	- Output ID if the image is any of the K persons

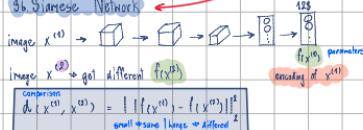
35 One-shot learning → learning from one example to recognize input image person → CNN → O (4 or 8 outputs)

learning a similarity param.

During, $\text{img}_1, \text{img}_2 = \text{degree of difference between images}$

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

36 Siamese Networks



37 Triplet Loss

Learning Objective: Anchor(A) Positive (P) + different neg.

Anchor(A) Negative (N)

Goal: $d(A, P) < d(A, N)$

$$\|f(A) - f(P)\|^2 < \|f(A) - f(N)\|^2$$

idea: to make sure the NN won't get confused → no margin

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 < 0$$

loss function: 10 pictures for 1 person better accuracy

given 3 images A, P, N

$$L(A, P, N) = \max\left(\frac{\|f(A) - f(P)\|^2}{m}, \frac{\|f(A) - f(N)\|^2}{m}\right)$$

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

Choosing the triplets A, P, N

- During training, if A, P, N are chosen randomly

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

is easily satisfied

- Choose triplets that're hard to train on. $d(A, P) \approx d(A, N)$

38 Face Verification

Learning the similarity function



39 What is neural style transfer?



40 What are deep CNNs learning?

- pick a unit in layer 1: Find the nine images patches (repeat for other units)
- that maximize the unit's activation

Visualizing deep layers



41 Cost function (for build neural style transfer)

Content(C) + Style(S) \Rightarrow Generated image(G)

Content function $J(G) = \alpha J(C, G) + \beta J(S, G)$
Style function $J(G) = \alpha J(C, G) + \beta J(S, G)$

Find the generated image G

1 Initialize G randomly. $G: 100 \times 100 \times 3$

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

$$G = G - \frac{d}{dG} (J(G))$$

42 Content Cost Function

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

- say using hidden layer l to compute cost
- use pre-trained ConvNet (e.g. vgg network)
- let $a^{F(l)}$ and $a^{G(l)}$ be the activation of layer l
- If $a^{F(l)}$ and $a^{G(l)}$ are similar then they will have

$$J_{content}(C, G) = \|a^{F(l)} - a^{G(l)}\|_F^2$$

43 Style Cost Function