

Computer Vision

Convolution NN

Rapid growth

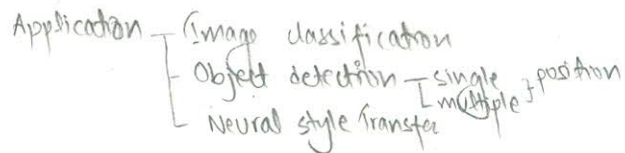
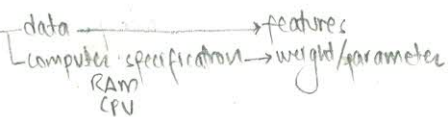
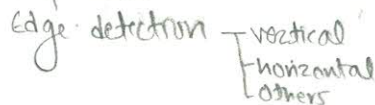
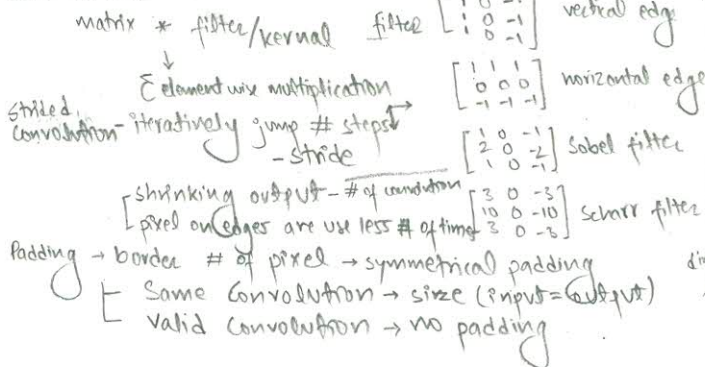


Image Classification



Convolution



Shades

- light to dark detect - 1 (absol)
- dark to light detect - 0

learning parameter

central position

- natural dimension

Convolution → Cross-correlation

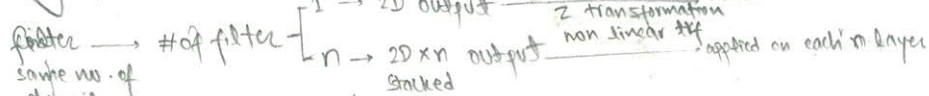
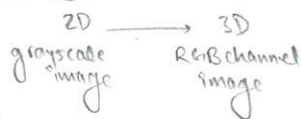
sliding dot product

sliding inner product

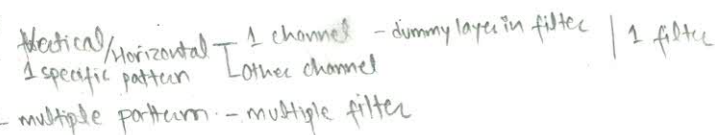
+ diagonal flip of filter → convolution

double mirroring operation - associative property - Signal processing

Convolutions over volume



Edge detector



Type of Layers in Convolution NN



Pooling Layer

Max Pooling

Average Pooling

feature detected over region gets max

average

3D

apply on each channel independently

filter size

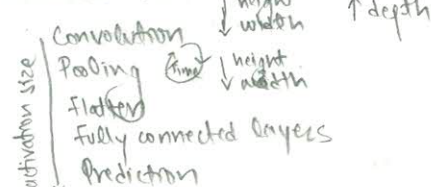
stride

no learning while backward prop

data → efficient

robust feature

Convolution NN



Why convolution

no. of parameter ↓

overfitting

Parameter sharing - filter

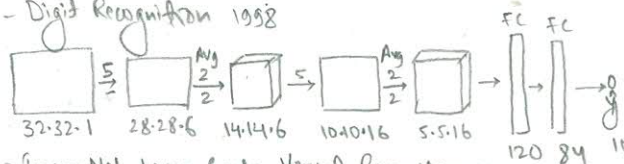
Spatially of connection - filter

CASE STUDIES

NN Architecture

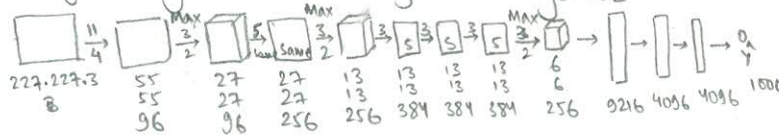
Classic Networks

LeNet5 - Digit Recognition 1998



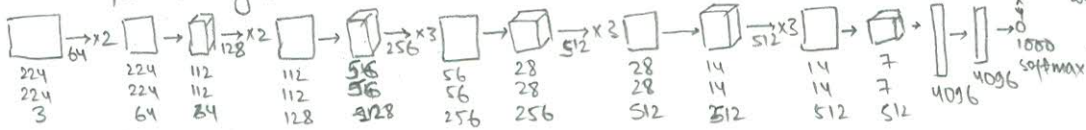
60,000 parameters
Pooling → average
Activation function → Sigmoid, tanh
Classifier → softmax
Valid Convolution

AlexNet - ImageNet Large Scale Visual Recognition Challenge - 2012



60 million params
ImageNet dataset
ReLU
Local Response Normalization

VGG-16 16 layers with weights 2015

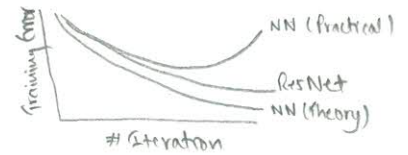
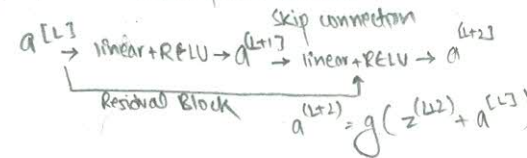


Convolution filter 3×3
stride = 1, same
Max Pool 2×2 , $s=2$
138 million params
Easy architecture
- filter

ResNets Networks 2015

Deep NN - Vanishing/Exploding Gradients

Plain Network main path + shortcut



Parameter - Easy to learn Identity function
- best performance

Condition - same dimension - same convolut
- same dimension - weight transfer

Inception Networks 2014

1x1 Convolution Network in Network

choose - convolution - layer 2
- pooling

> 1 channel → reduce channel - filters
non-trivial + non linearity
can & computational cost

Stack - computational cost - 1x1 convolution - bottle neck layer
Channel concat

side branch → prediction
- overfitting
- regularization
- hidden layer → prediction

Transfer Learning → Opensource Implementation - github

Codes & weight
Architecture

→ use change

small data - softmax function
medium data - last few layers - change re tune
large data - full network - re train

freeze trainable parameters
transform input → save

Data Augmentation

Shape & size - mirroring
- random cropping
- rotation
- shearing
- local wrapping
Color shifting - PCA color augmentation

Implementation
CPU thread

State of Computer Vision

less data - object detection
more data - image recognition
- simple Algorithm
- hand engineering features

Source of knowledge
- labelled data
- hand engineering

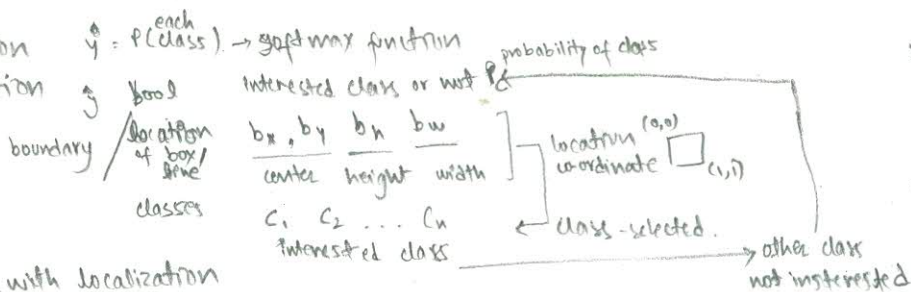
Benchmark

Ensembling
Multi step at test time

Deployment
Production
computation cost

OBJECT DETECTION

Image classification
Landmark Detection
(Snap chat)



loss function

$$L(\hat{y}, y) = \begin{cases} \sum_{i=1}^n (\hat{y}_i - y_i)^2 & \text{if object} = \text{True} \\ 0 & \text{all} \\ (\hat{y}_i - y_i) & \text{if object} = \text{False} \end{cases}$$

only bool

Object detection with localization

Classified Images \rightarrow search \rightarrow canvas

Sliding windows detection algorithms 2014

window

stride/step size

size - multiple sequentially

Simple Classifier \rightarrow linear regression + hand engineered features

computation cost \rightarrow CNN

Convolutional Implementation - share computation

CNN fully connected layer \rightarrow convolution layer

$1 \times 1 \times n_c \times n_f$ convolution multiple layers

convolution - windows pooling - stride

Object Detection

YOLO (You look only once) algorithm

Image \rightarrow grid cells \rightarrow each - target label

P_c - one object \rightarrow one grid center of object

b_x, b_y, b_h, b_w relative to grid cell

Bounding box accuracy \downarrow

position \downarrow stride

Shape

- Boundary box of any aspect ratios explicitly incorporate

- convolutional implementation - real time

Multiple detection of same object

clean up detection - Non Max Suppression

- Select boxes (with off of P_c)
- Pick box with largest P_c
- Discard boxes with high similarity

Multiple Objects - Anchor Box

predefined shapes \rightarrow stack in target label

each object (training) - grid cell - mid point

Multiple obj in anchor boxes - tie breaker

Different obj - same grid cells - Anchor box

box shape - predefined hand engineered

K-means to detect obj shape

Region proposals - ignore non interesting regions

Subset few windows \rightarrow run classifier

segmentation algorithms - find blob - run classifier

R-CNN Region-label-boundary box

First R-CNN

Faster R-CNN

slow

convolutional implementation

Box shape - scale

Face Recognition + (Liveness)

Image \rightarrow Database
 \downarrow
Generate ID

Face Verification
Image \rightarrow ID
 \downarrow
Check

Learn from one example - One shot learning problem
CNN \rightarrow accuracy \downarrow scalable + extra person \rightarrow update model
Learning similarity function
- degree of differences \rightarrow cutoff

Siamese Network Architecture - 2014 DeepFace

Image $x^{(i)}$ $\xrightarrow{\text{CNN}}$ Vector $F(x^{(i)})$
encoding
representation of image

Objective 2015 - FaceNet
Triplet Image \rightarrow Anchor (A) Positive (P) Negative (N)
 \rightarrow same person \rightarrow different person

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

$d(x^{(i)}, x^{(j)}) = \|f(x^{(i)}) - f(x^{(j)})\|^2$
same param from CNN/encoding \rightarrow find encoding
- degree of diff - small \rightarrow same person
- degree of diff - large \rightarrow diff person

Loss function \rightarrow Triplet

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

$$J = \sum L(A, P, N)$$

choose triplet image \rightarrow hard to train - CNN
 \rightarrow Random

Binary Classification Problem

Image $\xrightarrow{\text{CNN}}$ Embedding \rightarrow logistic regression

- encoding
- difference in encoding - absolute diff
- x^2 similarity

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

Neural Style Transfer

Content (C) style (S)

Generate Image (G)
random assign
keep updating

Cost function

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

Generate Image
 $G = G - \frac{\partial}{\partial C} J(G)$

Visualization of NN

Image \rightarrow Patches \rightarrow Activation Unit \rightarrow Repeat
given layer
Find patches with high activation \rightarrow visualize
max

Content Cost function - apply pre-trained ConvNet \rightarrow select layer
obj \rightarrow C & G - similar $J_C = \|a^{[l]}(C) - a^{[l]}(G)\|^2$

Style Cost function - correlation across channels
Style matrix $a^{[l]}(S)$
activation height width \rightarrow layer \rightarrow channel

given layer - $J_S = \|G^{[l]}(S)\|^2$
multiple layers - $J_S = \sum_l \lambda^{[l]} J_S^{[l]}(S, G)$
hyperparameter

shallow layer - simple feature
edges colors
deep layer - complex pattern

gram matrix - G

$$G_{K \times K} = \sum_i \sum_j a_{ijk} a_{ijl}$$

$$J_S = \|G^{[l]}(S) - G^{[l]}(G)\|^2$$

convolution - Data

1D - ECG Report 1D
2D - picture 2D
3D - CAT scan Movie with Time 3D

- channel

convolution