

Deep Learning for Computer Vision

Andrii Liubonko

Samsung R&D Institute Ukraine

July 24, 2018

Course overview

Unit I *DL overview. Classification. Modern CNNs.*

Unit II *Object Localization. Detection. Semantic Segmentation.*

Unit III *Metrics learning. Perspective.*

Contents

Deep Learning for Computer Vision.

Quick and dirty ML overview.

Basic elements of CNN

Motivation

Structure

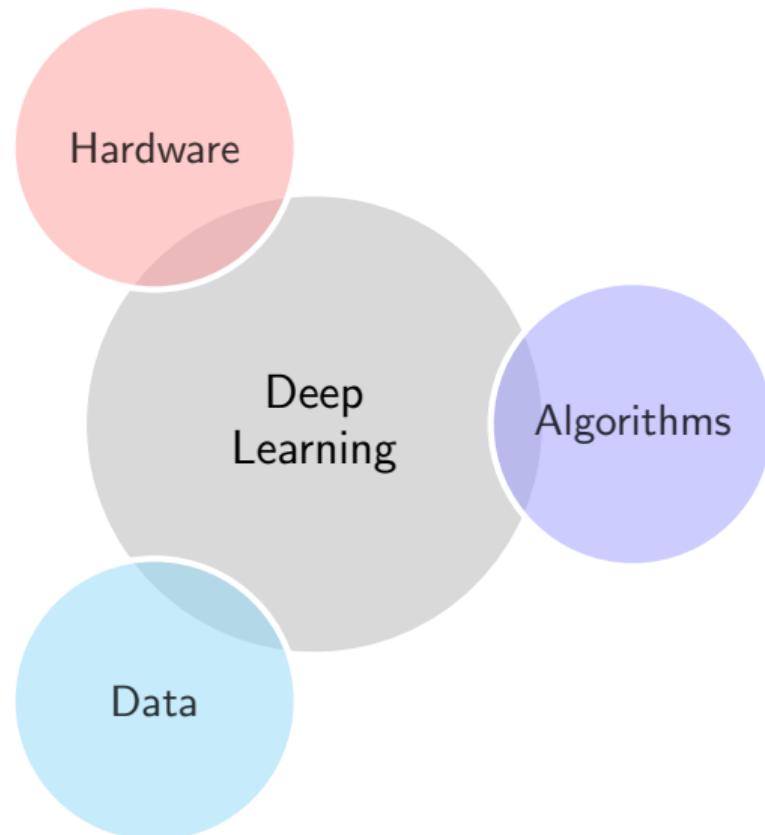
Structure [convolution layer]

Structure [fully-connected layer]

Case study

Summary

Deep Learning for Computer Vision.



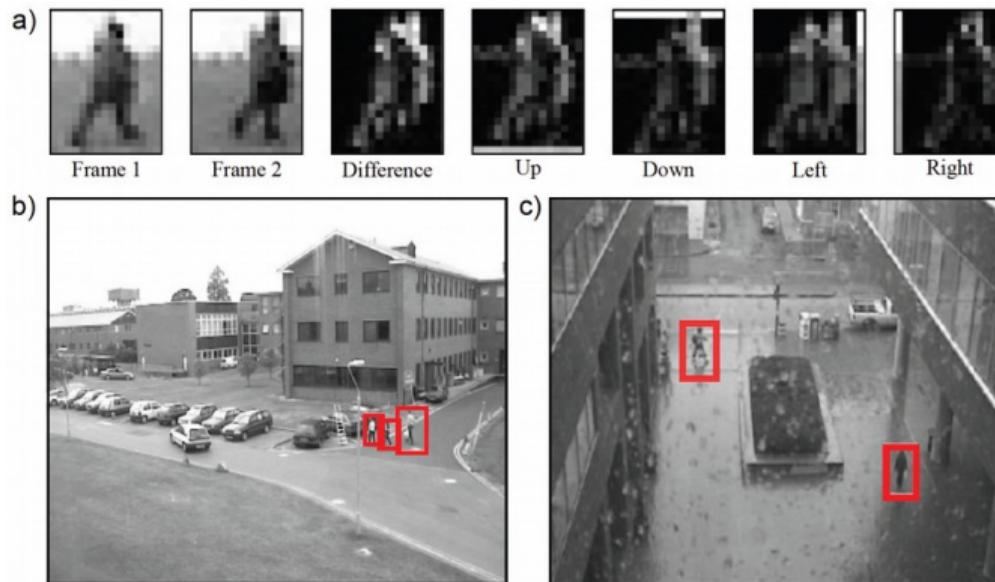
Deep Learning for Computer Vision.

Classification [e.g genre classification]



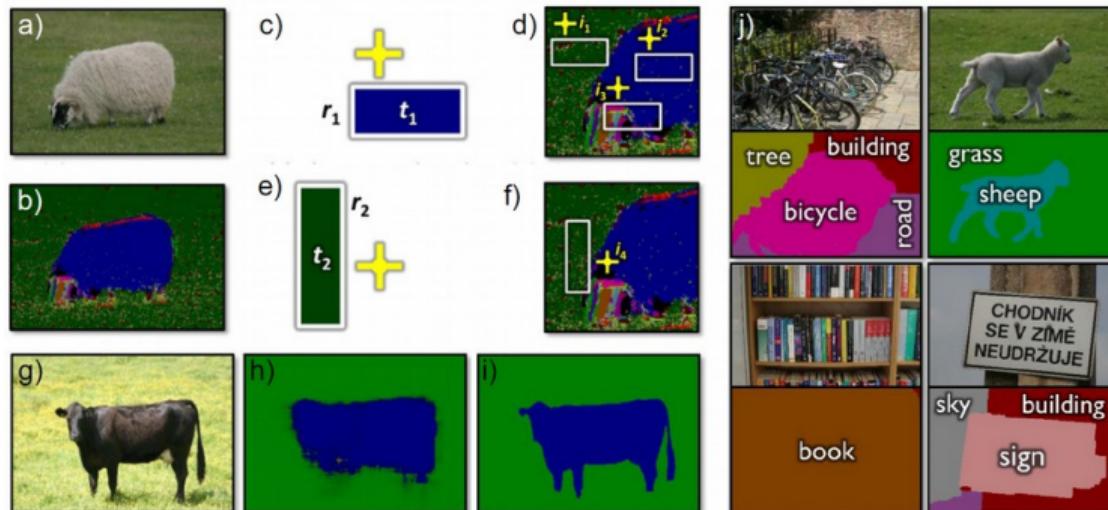
Deep Learning for Computer Vision.

Object detection [e.g. pedestrian detection]



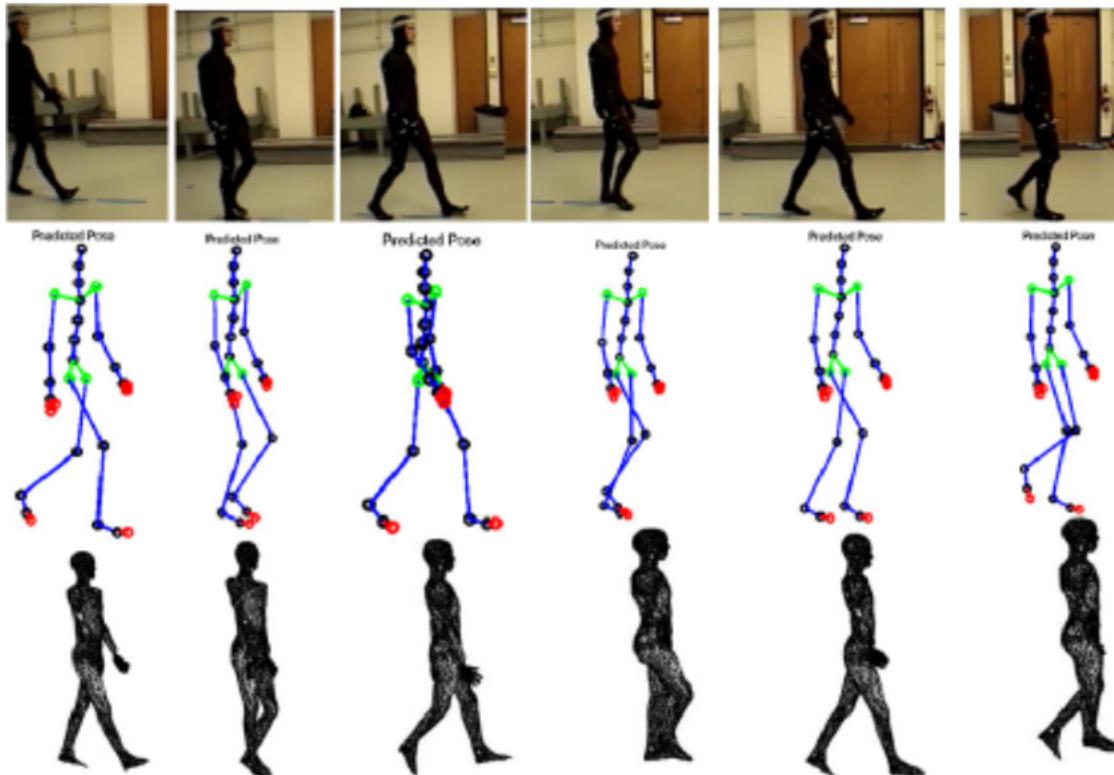
Deep Learning for Computer Vision.

Semantic segmentation



Deep Learning for Computer Vision.

Human pose estimation



Contents

Deep Learning for Computer Vision.

Quick and dirty ML overview.

Basic elements of CNN

Motivation

Structure

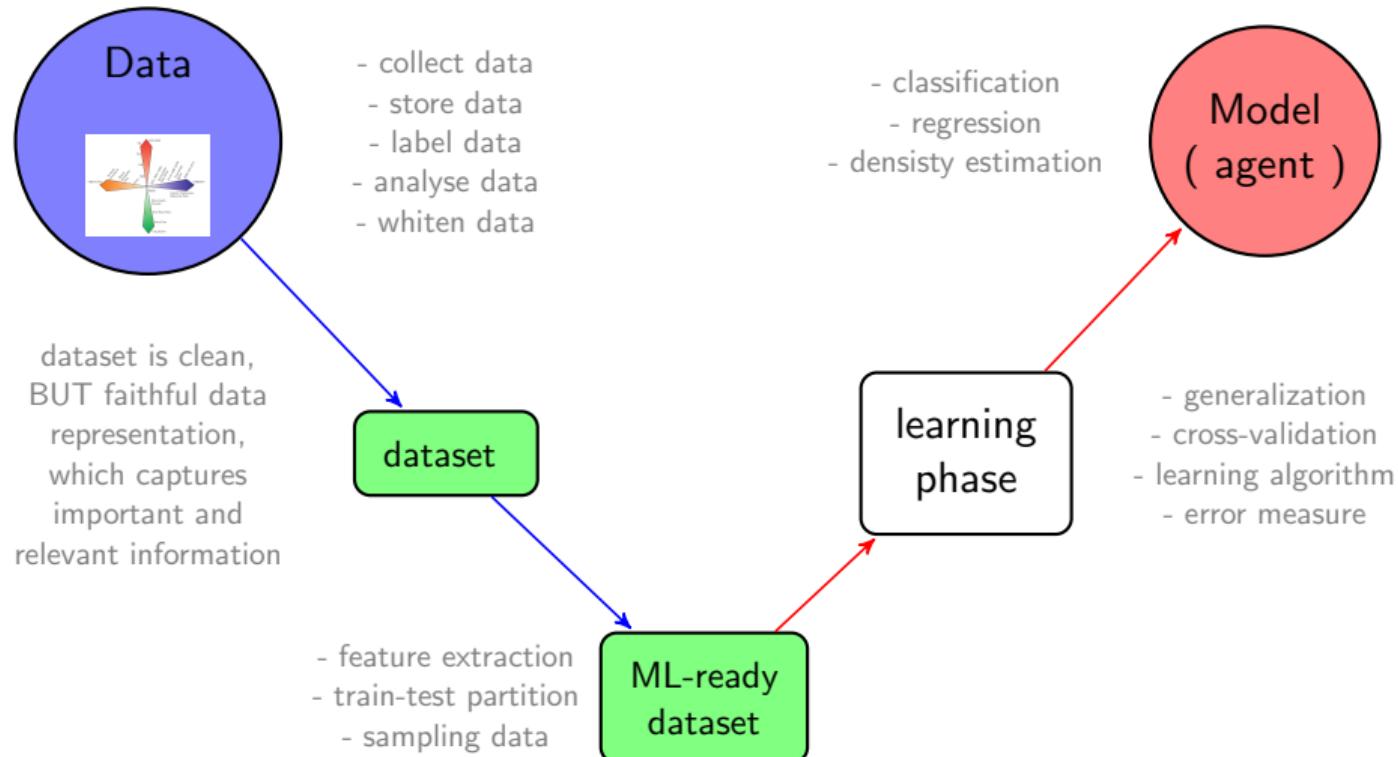
Structure [convolution layer]

Structure [fully-connected layer]

Case study

Summary

Quick and dirty ML overview.



Quick and dirty ML overview.[Types of learning]

- ▶ Supervised
- ▶ Unsupervised
- ▶ Semi-supervised
- ▶ RL

Quick and dirty ML overview.[Types of learning]

The goal of **supervised learning** is to learn a mapping from inputs \mathbf{x} to outputs y , given a labeled set of input-output pairs

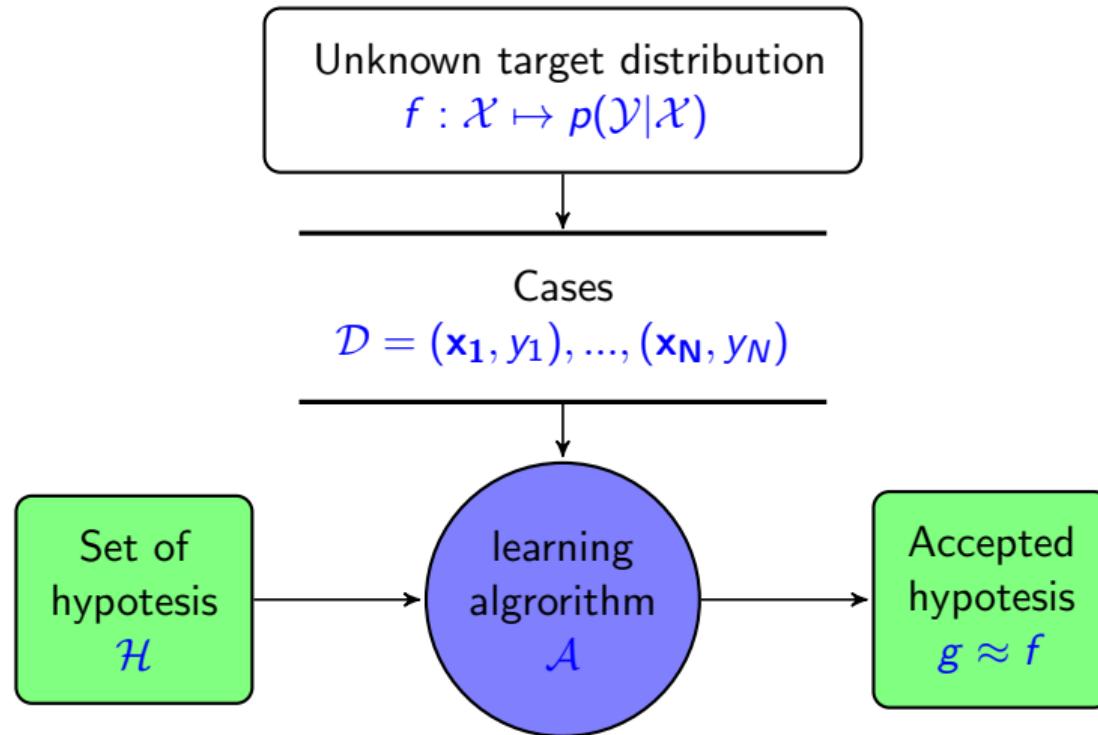
$$\mathcal{D} = (\mathbf{x}_i, y_i)_{i=1}^N$$

\mathcal{D} is called the **training set**, and N is the number of training examples.

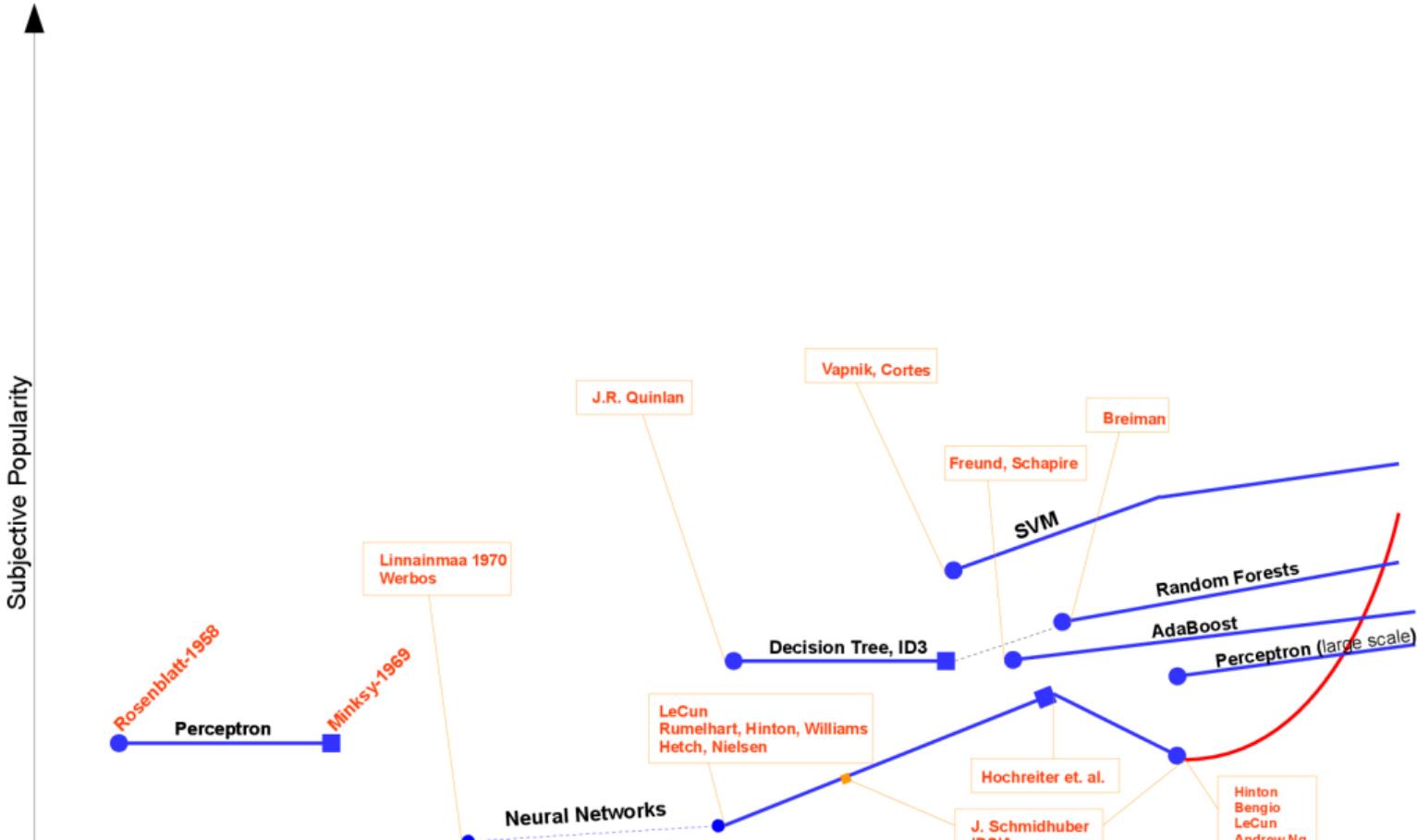
In the simplest setting, each training input \mathbf{x}_i is a D -dimensional vector of numbers. These are called **features** or attributes.

In general \mathbf{x}_i could be a complex structured object, such as an image, a sentence, an email message, a time series, a molecular shape, a graph, etc.

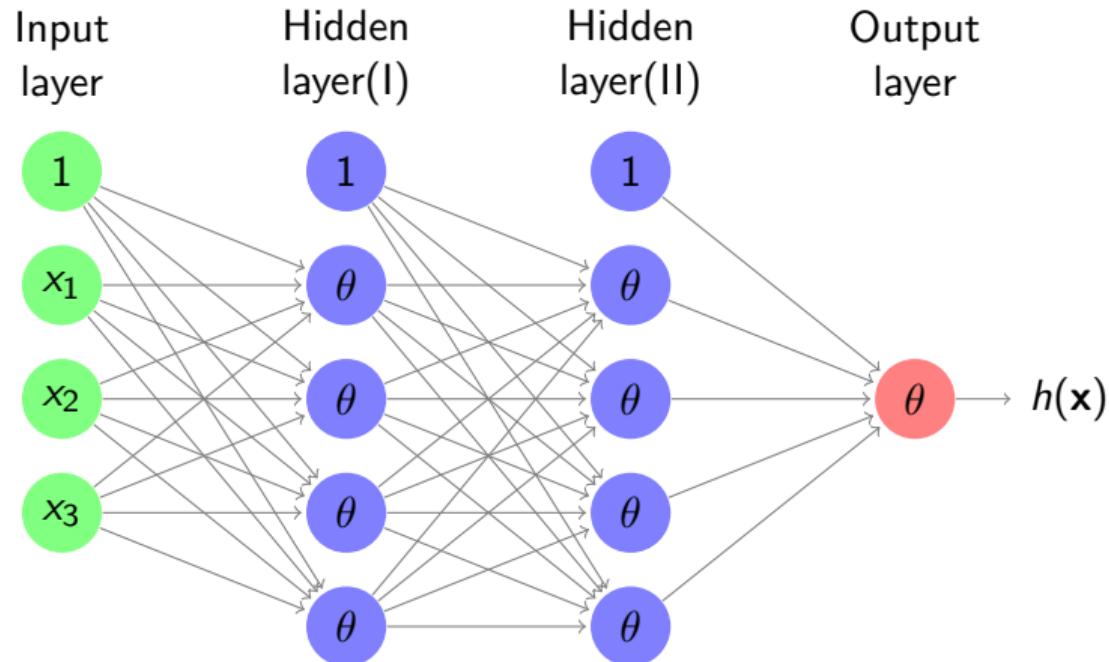
Quick and dirty ML overview.[Types of learning]



Quick and dirty ML overview.[Types of learning]



Quick and dirty ML overview. [NN - Architecture]



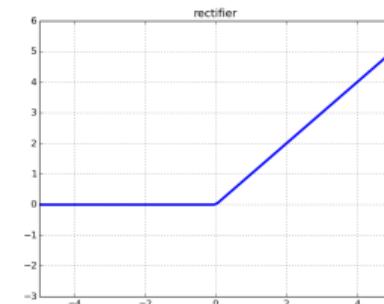
Quick and dirty ML overview. [NN - Architecture]

parameters in NN:

$$\omega_{ij}^{(l)} = \begin{cases} 1 \leq l \leq L & \text{layers} \\ 0 \leq i \leq d^{(l-1)} & \text{inputs} \\ 1 \leq j \leq d^{(l)} & \text{outputs} \end{cases}$$

activation:

$$x_j^{(l)} = \theta(s_j^{(l)}) = \theta\left(\sum_{i=0}^{d^{(l-1)}} \omega_{ij}^{(l)} x_i^{(l-1)}\right)$$



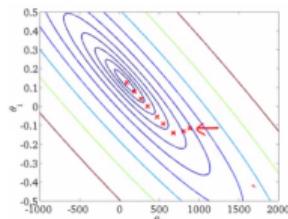
$$\theta(s) = \text{ReLU}(s) = \max(0, x)$$

Quick and dirty ML overview. [NN - Architecture]

Define **Loss (or Cost) function:**

$$L_{logloss} = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^M y_{nk} \cdot \log(p_{nk})$$

$$L_{rmse} = \frac{1}{N} \sum_{n=1}^N \| F(\mathbf{x}_n) - y_i \|_2$$



Gradient Descent (GD) minimizes:

$$L_{train}(\omega) = \frac{1}{N} \sum_{n=1}^N e(F(\mathbf{x}_n), y_n)$$

by iterative steps along $-\nabla L_{train}$:

$$\Delta\omega = -\eta \nabla L_{train}(\omega)$$

$$\omega_{prev} = \omega_{next} + \Delta\omega$$

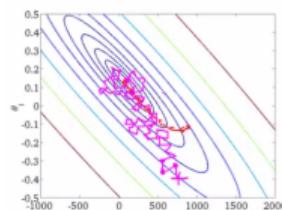
If $\nabla L_{train}(\omega)$ is based on all examples $\{\mathbf{x}_n, y_n\}$ then it is called **batch gradient descent**

Quick and dirty ML overview. [NN - Architecture]

Define **Loss (or Cost) function:**

$$L_{logloss} = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^M y_{nk} \cdot \log(p_{nk})$$

$$L_{rmse} = \frac{1}{N} \sum_{n=1}^N \| F(\mathbf{x}_n) - y_n \|_2$$



Gradient Descent (GD) minimizes:

$$L_{train}(\omega) = \frac{1}{N} \sum_{n=1}^N e(F(\mathbf{x}_n), y_n)$$

by iterative steps along $-\nabla L_{train}$:

$$\Delta\omega = -\eta \nabla L_{train}(\omega)$$

$$\omega_{prev} = \omega_{next} + \Delta\omega$$

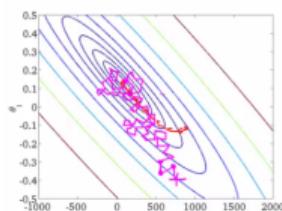
If $\nabla L_{train}(\omega)$ is based on all examples $\{\mathbf{x}_n, y_n\}$ then it is called
**stochastic (mini-batch)
gradient descent - SGD**

Quick and dirty ML overview. [NN - Architecture]

Define **Loss (or Cost) function:**

$$L_{logloss} = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^M y_{nk} \cdot \log(p_{nk})$$

$$L_{rmse} = \frac{1}{N} \sum_{n=1}^N \| F(\mathbf{x}_n) - y_n \|_2$$



Gradient Descent (GD) minimizes:

$$L_{train}(\omega) = \frac{1}{N} \sum_{n=1}^N e(F(\mathbf{x}_n), y_n)$$

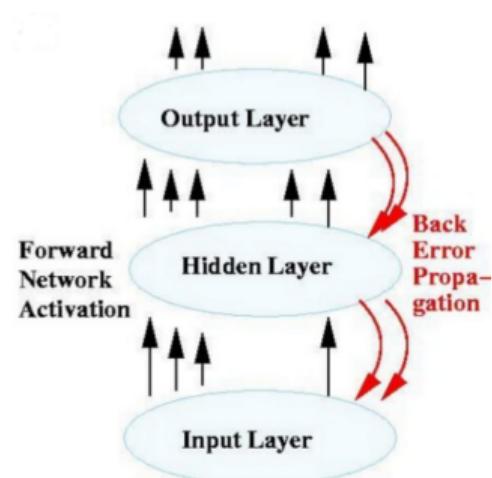
by iterative steps along $-\nabla L_{train}$:

$$\Delta\omega = -\eta \nabla L_{train}(\omega)$$

$$\omega_{prev} = \omega_{next} + \Delta\omega$$

If $\nabla L_{train}(\omega)$ is based on all examples $\{\mathbf{x}_n, y_n\}$ then it is called
momentum, rmsprop, adagrad, adam

Quick and dirty ML overview. [NN - Architecture]



Contents

Deep Learning for Computer Vision.

Quick and dirty ML overview.

Basic elements of CNN

Motivation

Structure

Structure [convolution layer]

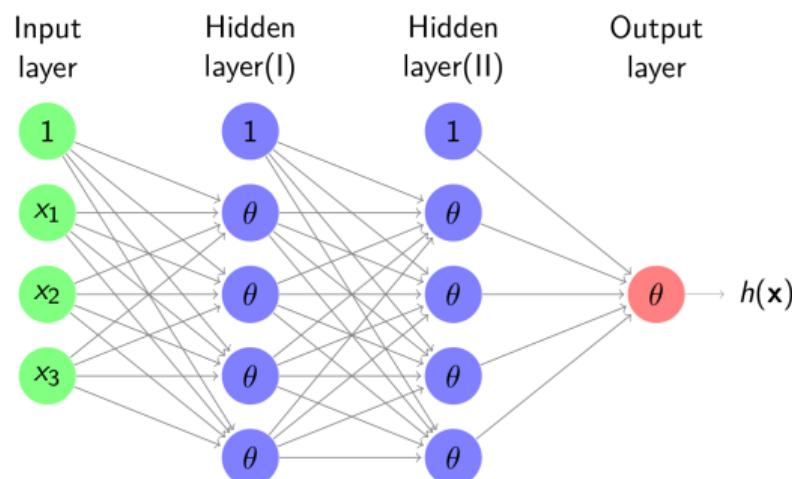
Structure [fully-connected layer]

Case study

Summary

Basic elements of CNN

Image processing with standard Feed-forward Neural Networks



Input layer
 $64 \times 64 = 4096$

Weights: $4096 \times$
(size of second
layer)

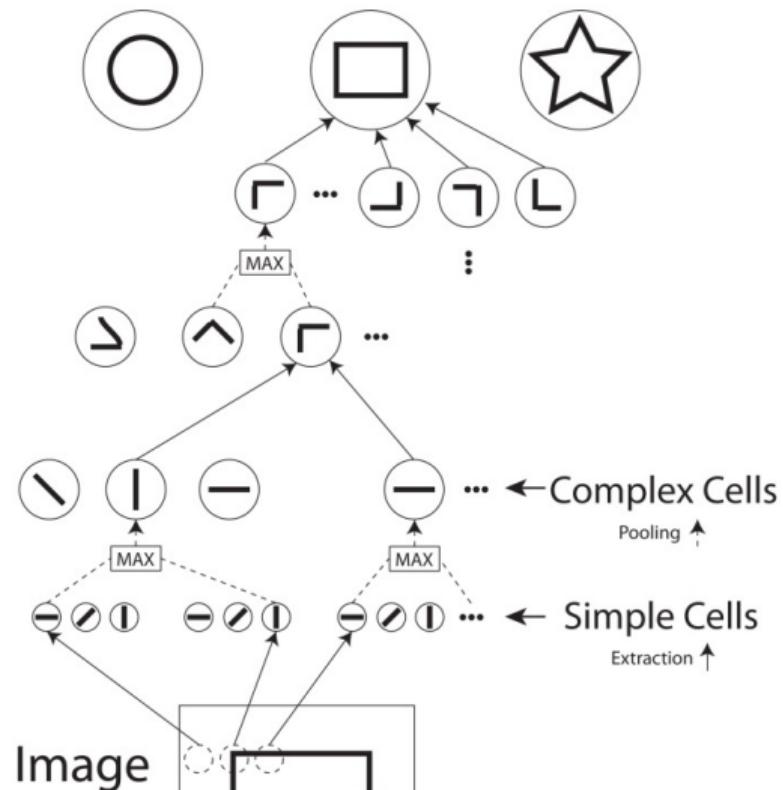
=> millions of
parameters (for
just one layer)

Need for
alternative
architecture

Basic elements of CNN

Motivation

Hubel & Wiesel (1960s) => Hierarchical organization



Basic elements of CNN

Motivation

Hubel & Wiesel (1960s) => Hierarchical organization



Contents

Deep Learning for Computer Vision.

Quick and dirty ML overview.

Basic elements of CNN

Motivation

Structure

Structure [convolution layer]

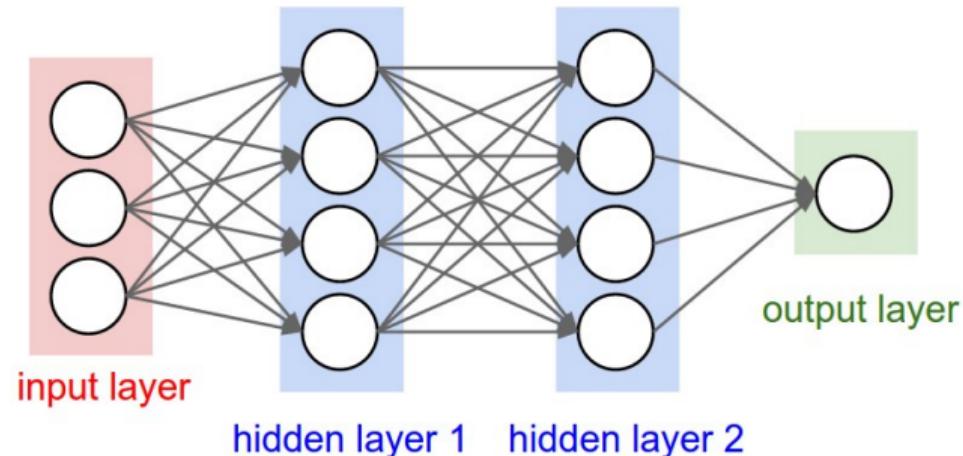
Structure [fully-connected layer]

Case study

Summary

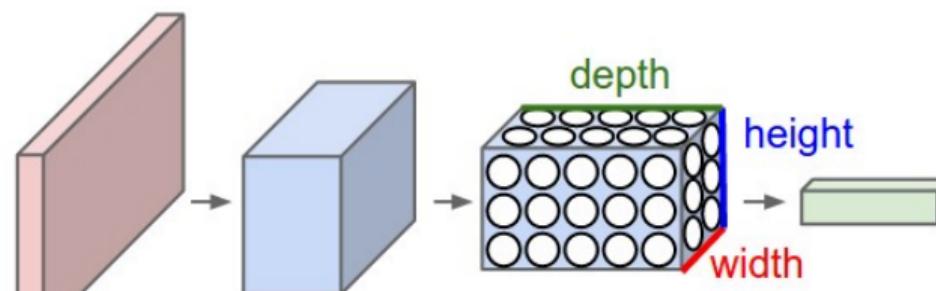
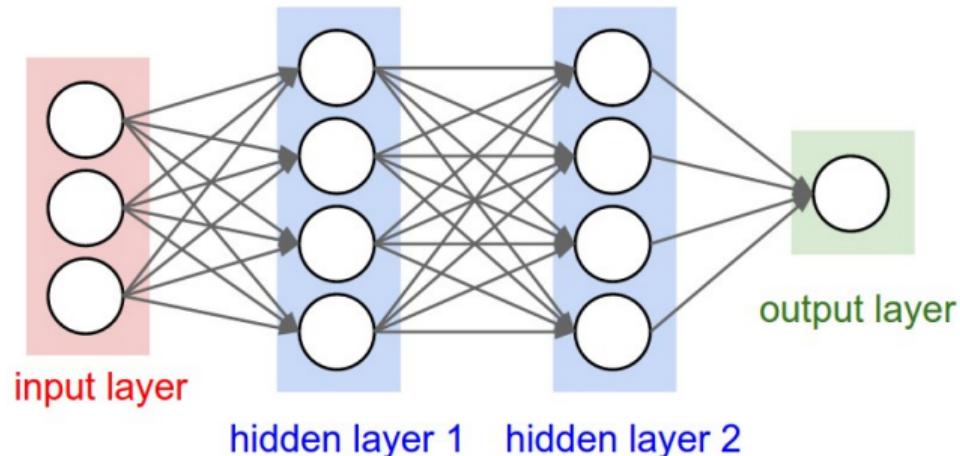
Basic elements of CNN

Structure



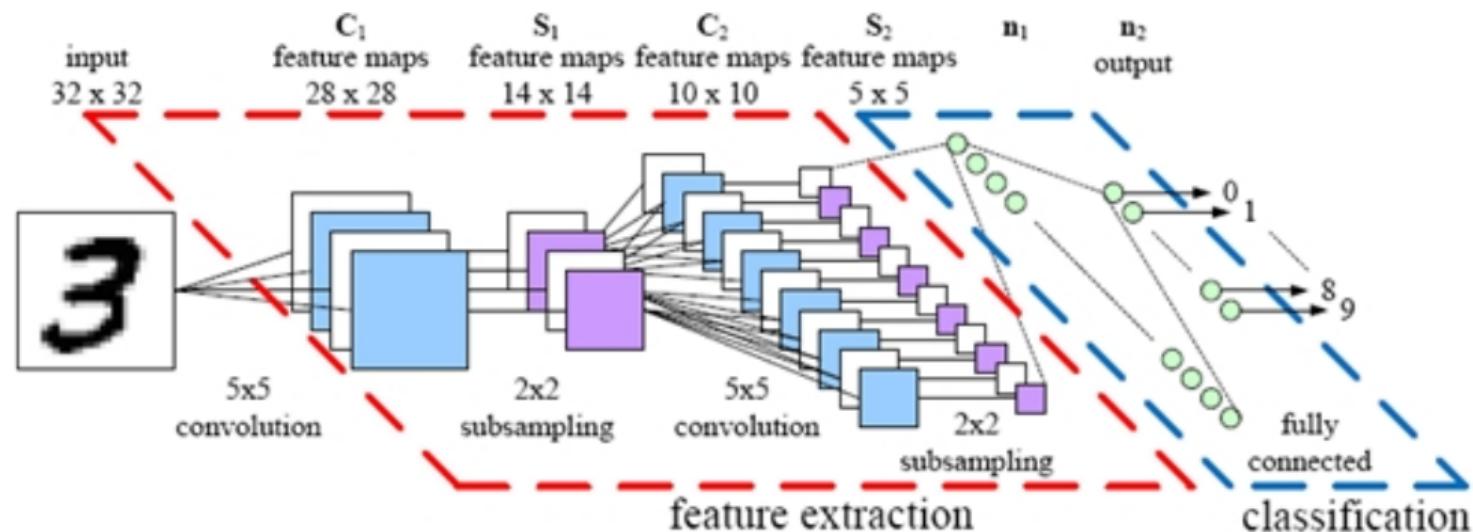
Basic elements of CNN

Structure



Basic elements of CNN

Structure



Basic elements of CNN

Structure

- ▶ INPUT holds the raw pixel values of the image.
- ▶ CONV layer computes the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and the region they are connected to in the input volume.
- ▶ POOL layer performs a downsampling operation along the spatial dimensions (width, height).
- ▶ FC (i.e. fully-connected) layer computes the class scores. As with ordinary Neural Networks and as the name implies, each neuron in this layer is connected to all the numbers in the previous volume.

Contents

Deep Learning for Computer Vision.

Quick and dirty ML overview.

Basic elements of CNN

Motivation

Structure

Structure [convolution layer]

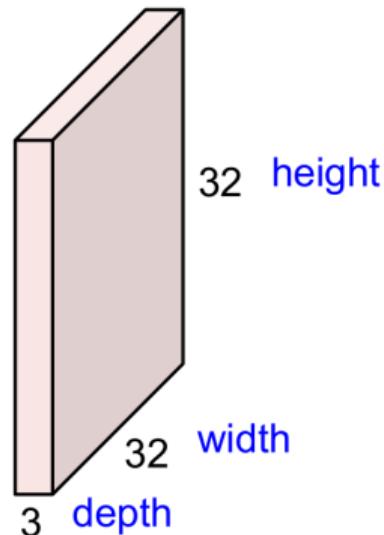
Structure [fully-connected layer]

Case study

Summary

Basic elements of CNN

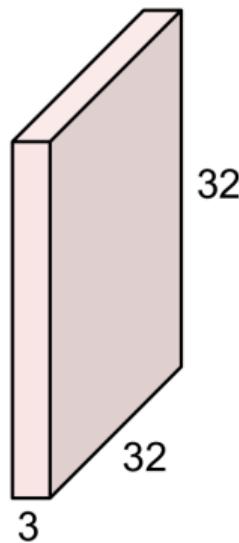
32x32x3 image



[from Fei-Fei Li & Andrej Karpathy & Justin Johnson "CNN for Computer Vision"]

Basic elements of CNN

32x32x3 image



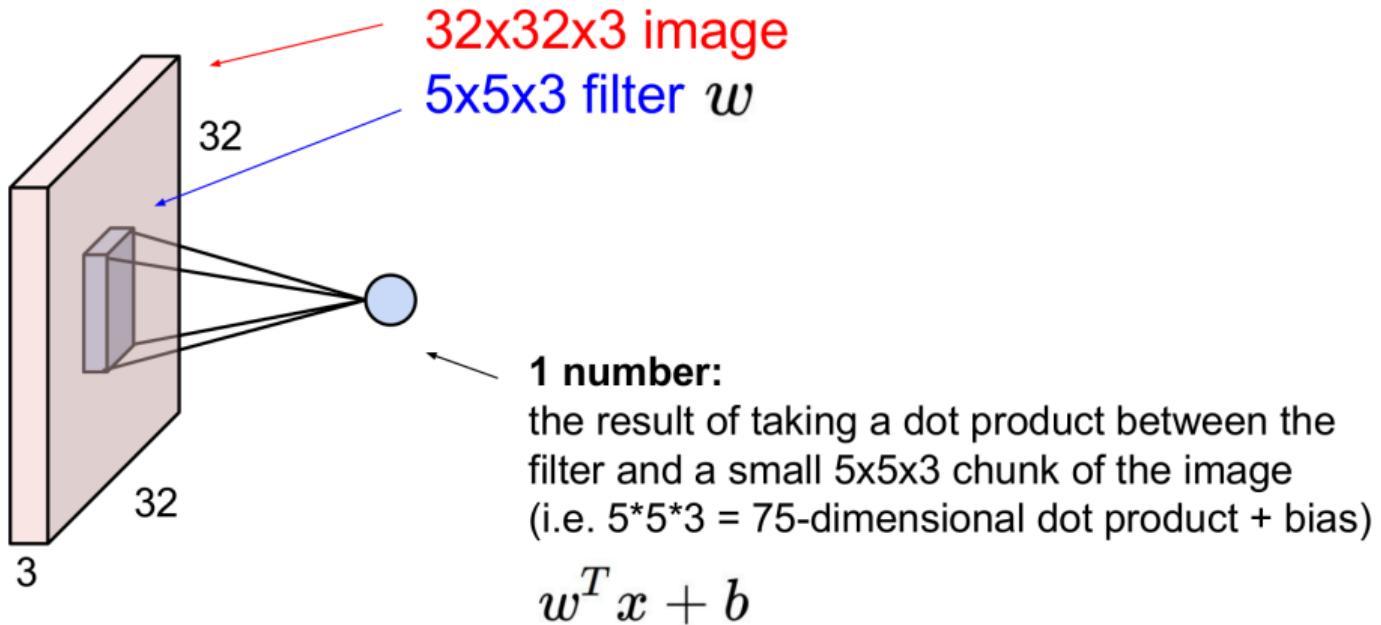
5x5x3 filter



Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

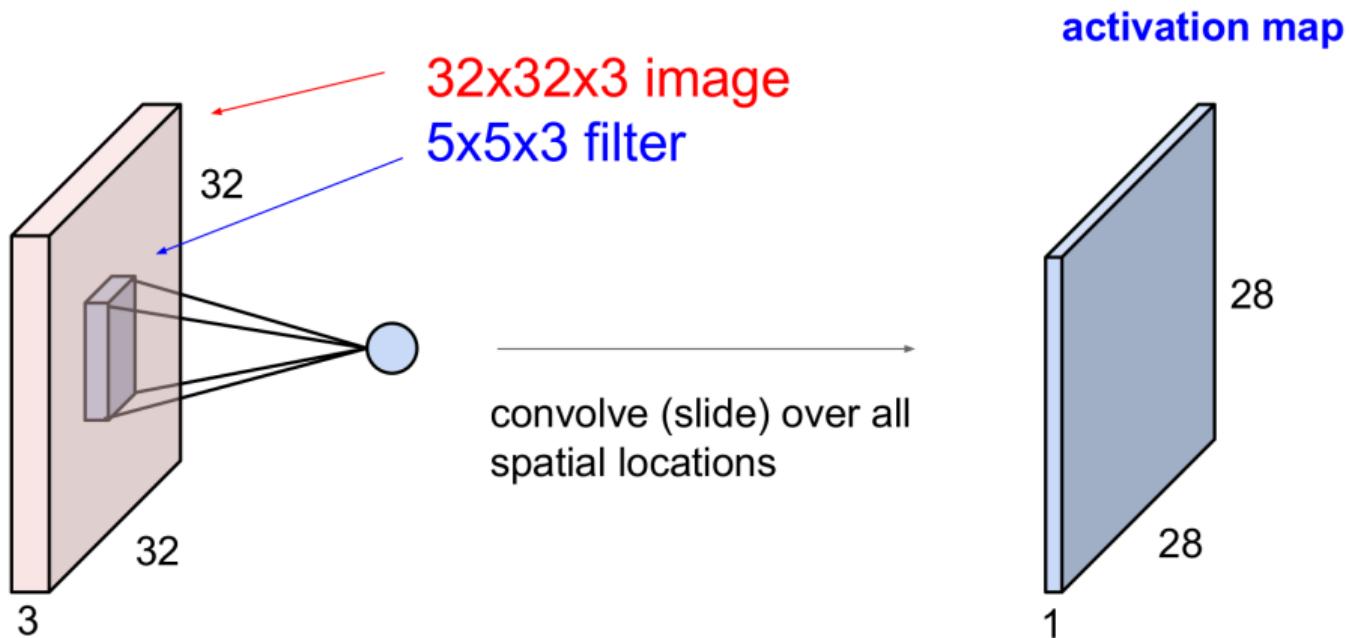
[from Fei-Fei Li & Andrej Karpathy & Justin Johnson "CNN for Computer Vision"]

Basic elements of CNN



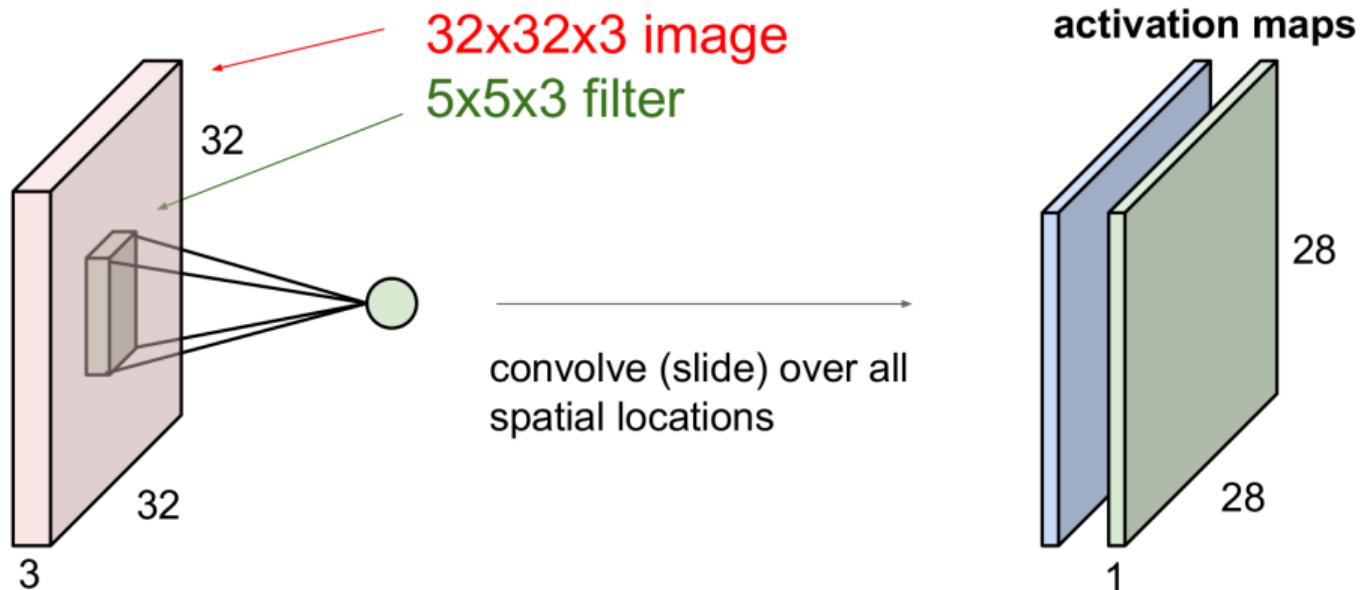
[from Fei-Fei Li & Andrej Karpathy & Justin Johnson "CNN for Computer Vision"]

Basic elements of CNN



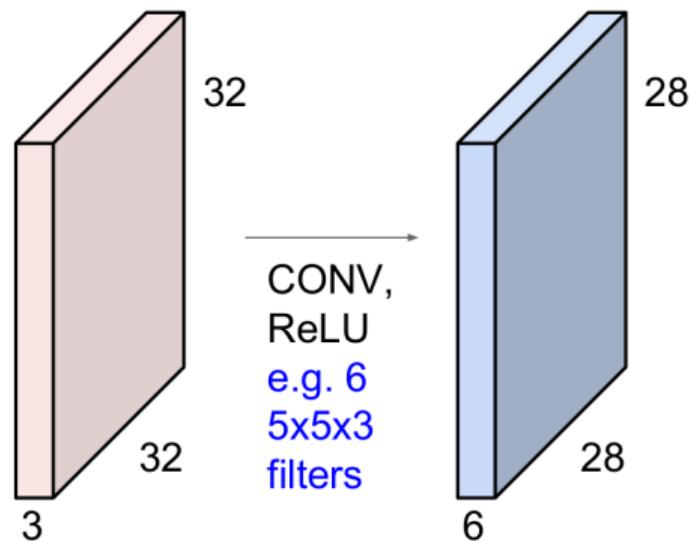
[from Fei-Fei Li & Andrej Karpathy & Justin Johnson "CNN for Computer Vision"]

Basic elements of CNN



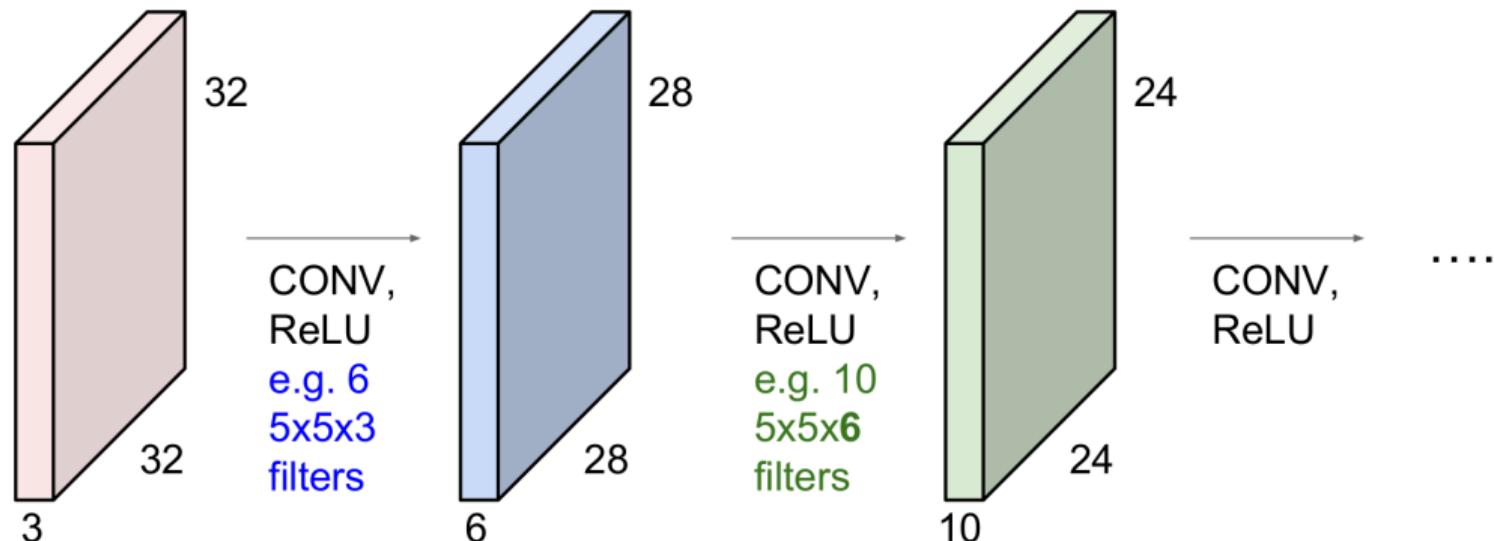
[from Fei-Fei Li & Andrej Karpathy & Justin Johnson "CNN for Computer Vision"]

Basic elements of CNN



[from Fei-Fei Li & Andrej Karpathy & Justin Johnson "CNN for Computer Vision"]

Basic elements of CNN

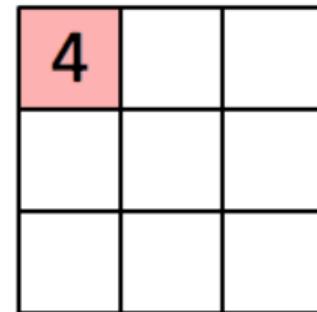


[from Fei-Fei Li & Andrej Karpathy & Justin Johnson "CNN for Computer Vision"]

Basic elements of CNN

1 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	0	0
0 <small>x0</small>	1 <small>x1</small>	1 <small>x0</small>	1	0
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image



Convolved
Feature

Basic elements of CNN

1	1 _{x1}	1 _{x0}	0 _{x1}	0
0	1 _{x0}	1 _{x1}	1 _{x0}	0
0	0 _{x1}	1 _{x0}	1 _{x1}	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	

Convolved
Feature

Basic elements of CNN

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	4

Convolved
Feature

Basic elements of CNN

1	1	1	0	0
0 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	1	0
0 <small>x0</small>	0 <small>x1</small>	1 <small>x0</small>	1	1
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	0
0	1	1	0	0

Image

4	3	4
2		

Convolved
Feature

Basic elements of CNN

1	1	1	0	0
0	1 <small>$\times 1$</small>	1 <small>$\times 0$</small>	1 <small>$\times 1$</small>	0
0	0 <small>$\times 0$</small>	1 <small>$\times 1$</small>	1 <small>$\times 0$</small>	1
0	0 <small>$\times 1$</small>	1 <small>$\times 0$</small>	1 <small>$\times 1$</small>	0
0	1	1	0	0

Image

4	3	4
2	4	

Convolved
Feature

Basic elements of CNN

1	1	1	0	0
0	1	1 _{x1}	1 _{x0}	0 _{x1}
0	0	1 _{x0}	1 _{x1}	1 _{x0}
0	0	1 _{x1}	1 _{x0}	0 _{x1}
0	1	1	0	0

Image

4	3	4
2	4	3

Convolved
Feature

Basic elements of CNN

1	1	1	0	0
0	1	1	1	0
0 x1	0 x0	1 x1	1	1
0 x0	0 x1	1 x0	1	0
0 x1	1 x0	1 x1	0	0

Image

4	3	4
2	4	3
2		

Convolved
Feature

Basic elements of CNN

1	1	1	0	0
0	1	1	1	0
0	0	1	1	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	4
2	4	3
2	3	

Convolved
Feature

Basic elements of CNN

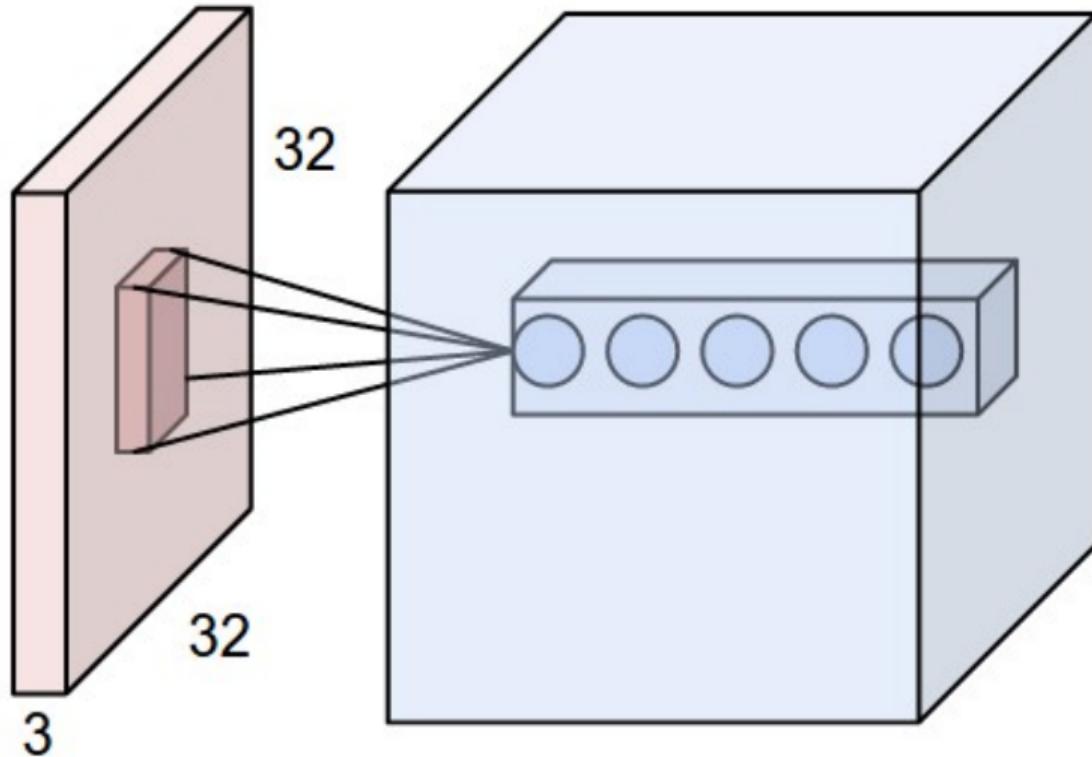
1	1	1	0	0
0	1	1	1	0
0	0	1 x1	1 x0	1 x1
0	0	1 x0	1 x1	0 x0
0	1	1 x1	0 x0	0 x1

Image

4	3	4
2	4	3
2	3	4

Convolved Feature

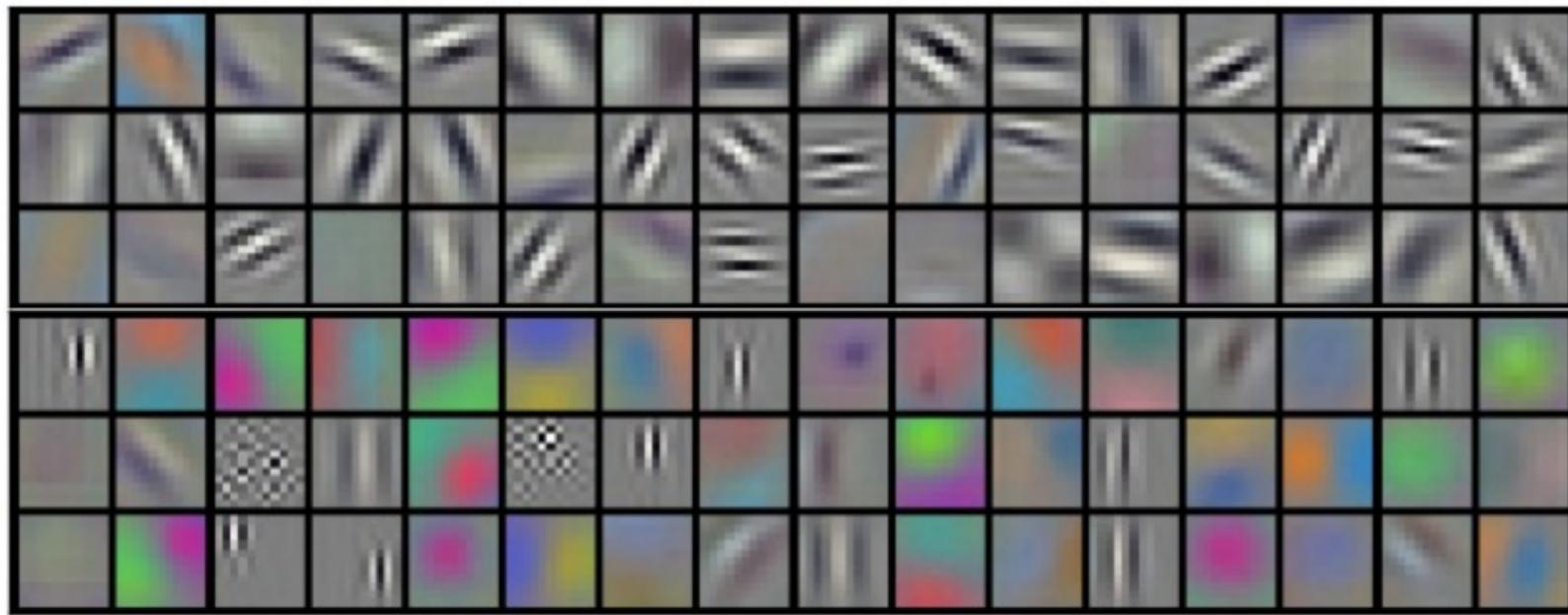
Basic elements of CNN



[from Fei-Fei Li & Andrej Karpathy & Justin Johnson "CNN for Computer Vision"]

Basic elements of CNN

Structure [convolution layer]



[from Fei-Fei Li & Andrej Karpathy & Justin Johnson "CNN for Computer Vision"]

Basic elements of CNN

- ▶ Accepts a volume of size $W1 \times H1 \times D1$
- ▶ Requires four hyperparameters:
 - ▶ Number of filters K ,
 - ▶ their spatial extent F ,
 - ▶ the stride S ,
 - ▶ the amount of zero padding P .
- ▶ Produces a volume of size $W2 \times H2 \times D2$ where:
 - ▶ $W2 = (W1 - F + 2P)/S + 1$,
 - ▶ $H2 = (H1 - F + 2P)/S + 1$
 - ▶ $D2 = K$
- ▶ With parameter sharing, it introduces $F \times F \times D1$ weights per filter, for a total of $(F \times F \times D1) \times K$ weights and K biases.
- ▶ In the output volume, the d -th depth slice (of size $W2 \times H2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

Contents

Deep Learning for Computer Vision.

Quick and dirty ML overview.

Basic elements of CNN

Motivation

Structure

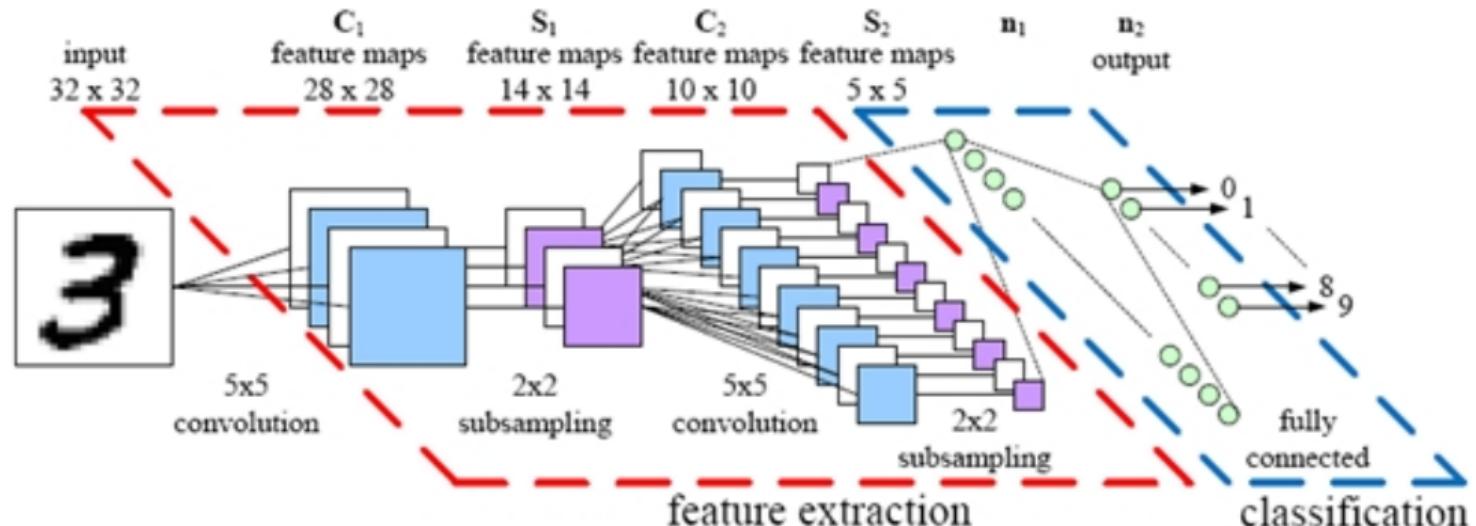
Structure [convolution layer]

Structure [fully-connected layer]

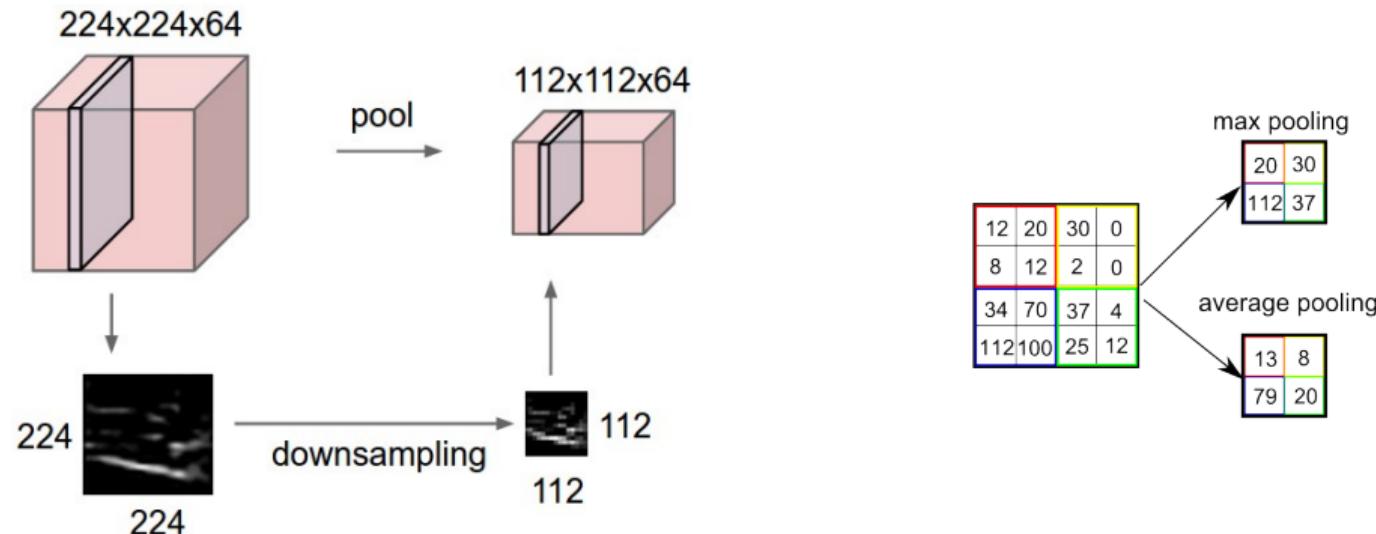
Case study

Summary

Basic elements of CNN



Basic elements of CNN



[from Fei-Fei Li & Andrej Karpathy & Justin Johnson "CNN for Computer Vision"]

Basic elements of CNN

- ▶ Accepts a volume of size $W1 \times H1 \times D1$
- ▶ Requires three hyperparameters:
 - ▶ their spatial extent F ,
 - ▶ the stride S ,
- ▶ Produces a volume of size $W2 \times H2 \times D2$ where:
 - ▶ $W2 = (W1 - F)/S + 1$
 - ▶ $H2 = (H1 - F)/S + 1$
 - ▶ $D2 = D1$

Contents

Deep Learning for Computer Vision.

Quick and dirty ML overview.

Basic elements of CNN

Motivation

Structure

Structure [convolution layer]

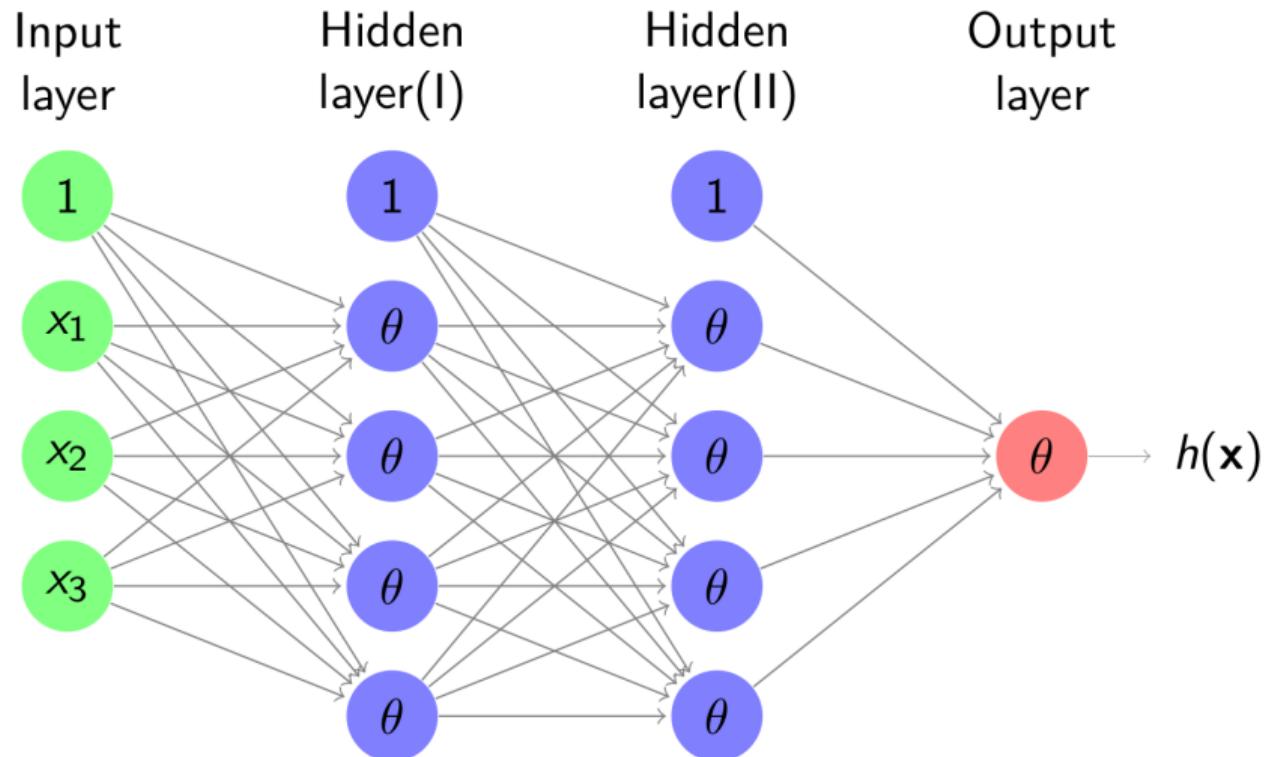
Structure [fully-connected layer]

Case study

Summary

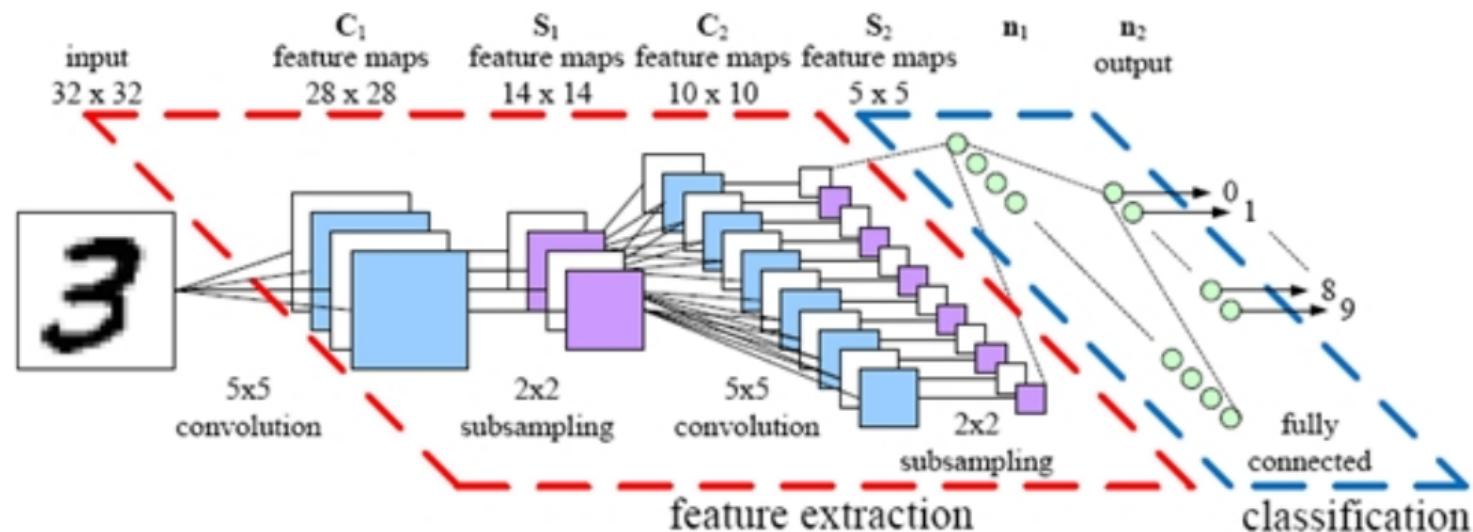
Basic elements of CNN

Structure [fully-connected layer]



Basic elements of CNN

Structure



Contents

Deep Learning for Computer Vision.

Quick and dirty ML overview.

Basic elements of CNN

Motivation

Structure

Structure [convolution layer]

Structure [fully-connected layer]

Case study

Summary

Basic elements of CNN

Case study

Case Study: AlexNet

[Krizhevsky et al. 2012]

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

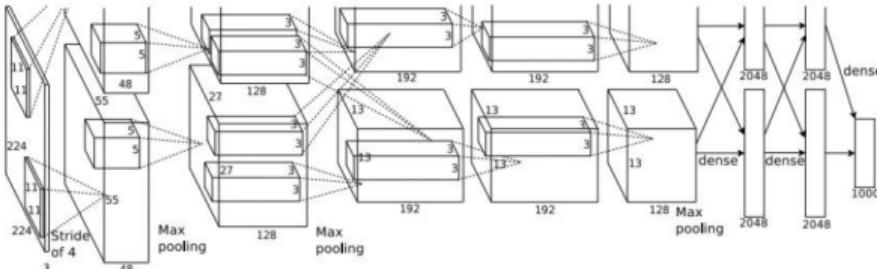
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% → 15.4%

Basic elements of CNN

Case study

INPUT: [224x224x3] memory: $224 \times 224 \times 3 = 150\text{K}$ params: 0 (not counting biases)
CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2\text{M}$ params: $(3 \times 3 \times 3) \times 64 = 1,728$
CONV3-64: [224x224x64] memory: $224 \times 224 \times 64 = 3.2\text{M}$ params: $(3 \times 3 \times 64) \times 64 = 36,864$
POOL2: [112x112x64] memory: $112 \times 112 \times 64 = 800\text{K}$ params: 0
CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6\text{M}$ params: $(3 \times 3 \times 64) \times 128 = 73,728$
CONV3-128: [112x112x128] memory: $112 \times 112 \times 128 = 1.6\text{M}$ params: $(3 \times 3 \times 128) \times 128 = 147,456$
POOL2: [56x56x128] memory: $56 \times 56 \times 128 = 400\text{K}$ params: 0
CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800\text{K}$ params: $(3 \times 3 \times 128) \times 256 = 294,912$
CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800\text{K}$ params: $(3 \times 3 \times 256) \times 256 = 589,824$
CONV3-256: [56x56x256] memory: $56 \times 56 \times 256 = 800\text{K}$ params: $(3 \times 3 \times 256) \times 256 = 589,824$
POOL2: [28x28x256] memory: $28 \times 28 \times 256 = 200\text{K}$ params: 0
CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400\text{K}$ params: $(3 \times 3 \times 256) \times 512 = 1,179,648$
CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
CONV3-512: [28x28x512] memory: $28 \times 28 \times 512 = 400\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
POOL2: [14x14x512] memory: $14 \times 14 \times 512 = 100\text{K}$ params: 0
CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
CONV3-512: [14x14x512] memory: $14 \times 14 \times 512 = 100\text{K}$ params: $(3 \times 3 \times 512) \times 512 = 2,359,296$
POOL2: [7x7x512] memory: $7 \times 7 \times 512 = 25\text{K}$ params: 0
FC: [1x1x4096] memory: 4096 params: $7 \times 7 \times 512 \times 4096 = 102,760,448$
FC: [1x1x4096] memory: 4096 params: $4096 \times 4096 = 16,777,216$
FC: [1x1x1000] memory: 1000 params: $4096 \times 1000 = 4,096,000$

TOTAL memory: $24\text{M} * 4 \text{ bytes} \approx 96\text{MB} / \text{image}$ (only forward! ~ 2 for bwd)

TOTAL params: 138M parameters

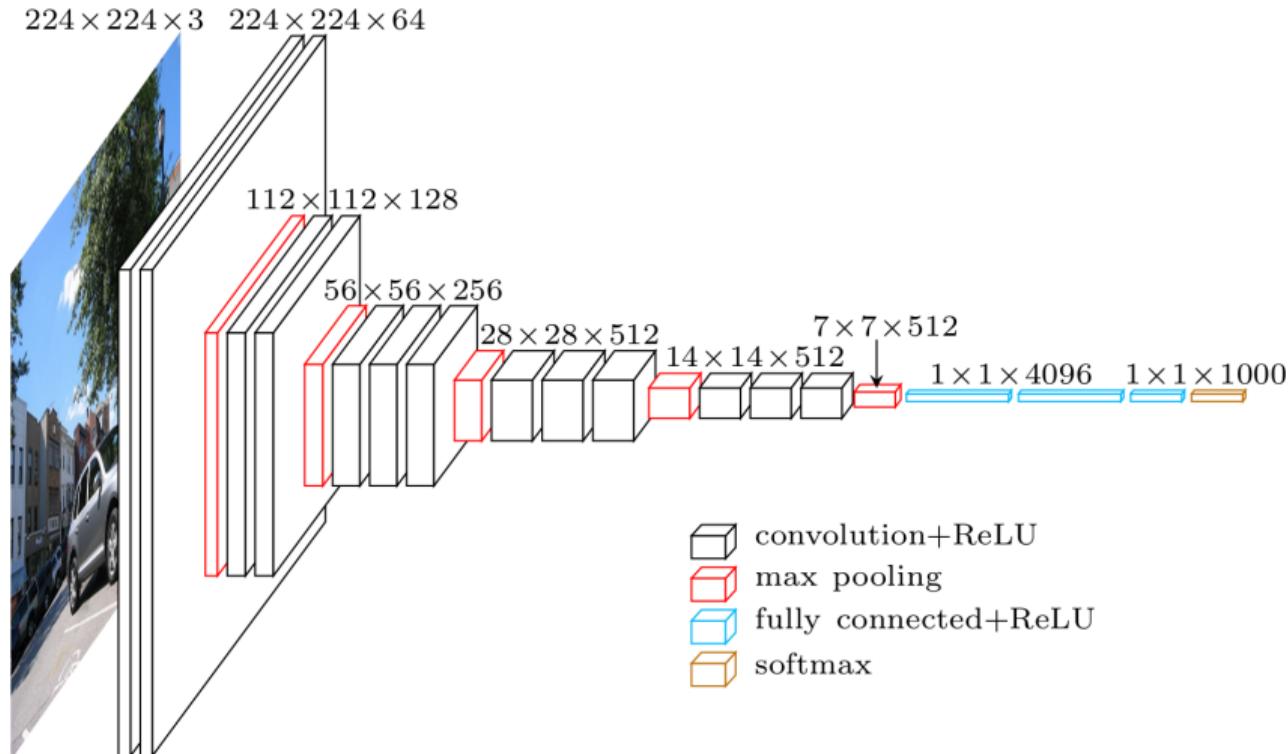


VGG16

Common names

Basic elements of CNN

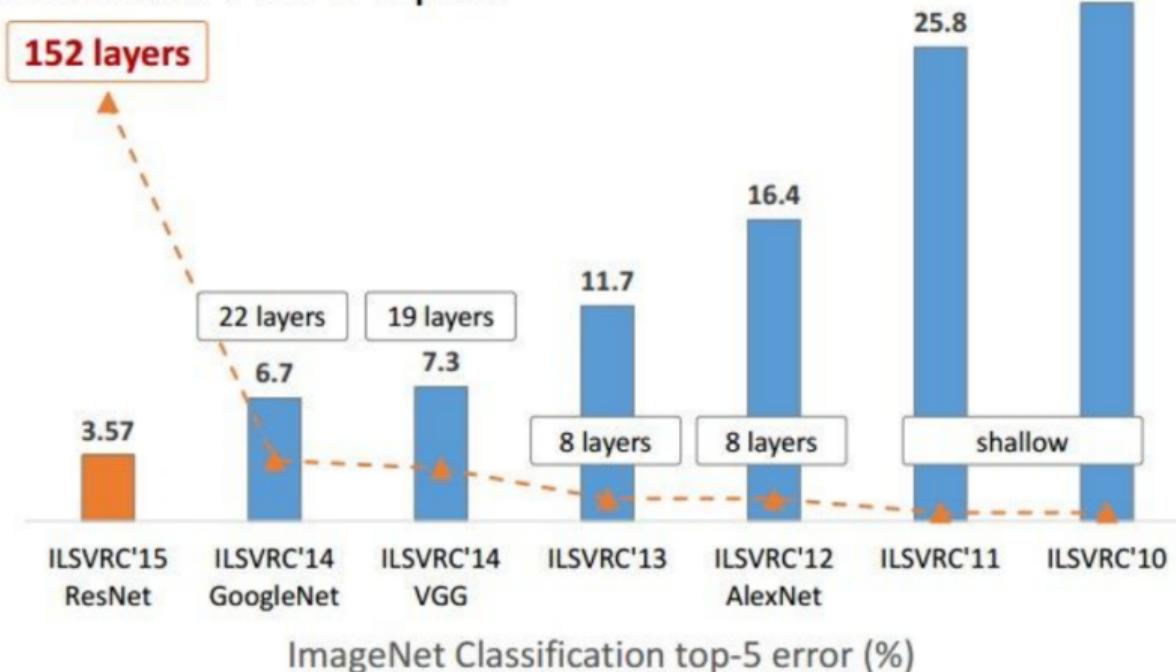
Case study



Basic elements of CNN



Revolution of Depth



Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



ResNet, 152 layers
(ILSVRC 2015)



Contents

Deep Learning for Computer Vision.

Quick and dirty ML overview.

Basic elements of CNN

Motivation

Structure

Structure [convolution layer]

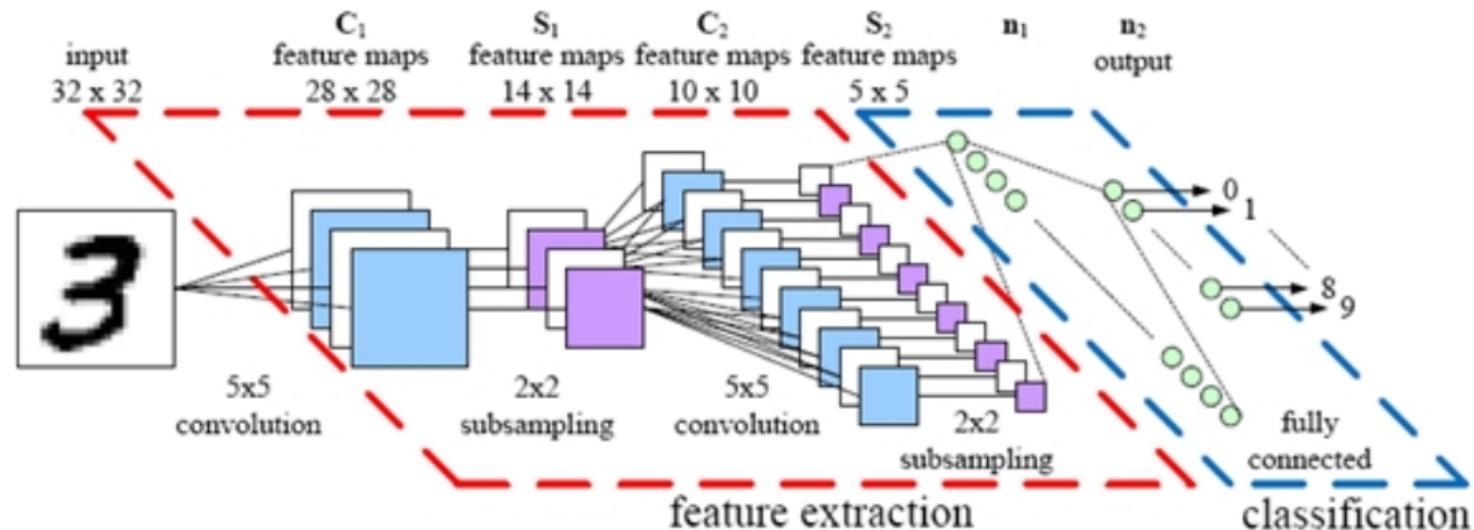
Structure [fully-connected layer]

Case study

Summary

Summary

CNN Architecture



Summary

CNN Architecture

- ▶ INPUT holds the raw pixel values of the image.
- ▶ CONV layer computes the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and the region they are connected to in the input volume.
- ▶ POOL layer performs a downsampling operation along the spatial dimensions (width, height).
- ▶ FC (i.e. fully-connected) layer computes the class scores. As with ordinary Neural Networks and as the name implies, each neuron in this layer is connected to all the numbers in the previous volume.