

# Computer Vision – Lecture 10

## Deep Learning

27.05.2019

Bastian Leibe  
 Visual Computing Institute  
 RWTH Aachen University  
<http://www.vision.rwth-aachen.de/>  
 leibe@vision.rwth-aachen.de

## Course Outline

- Image Processing Basics
- Segmentation & Grouping
- Object Recognition & Categorization
  - Sliding Window based Object Detection
- Local Features & Matching
  - Local Features – Detection and Description
  - Recognition with Local Features
- Deep Learning
- 3D Reconstruction

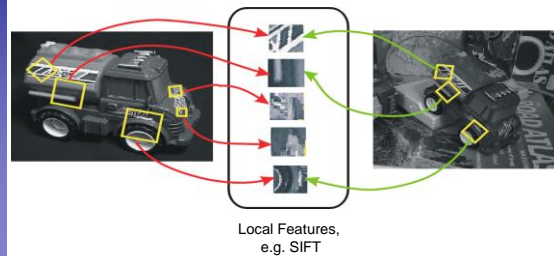
## Topics of This Lecture

- Recap: Recognition with Local Features
- Dealing with Outliers
  - RANSAC
  - Generalized Hough Transform
- Deep Learning
  - Motivation
  - Neural Networks
- Convolutional Neural Networks
  - Convolutional Layers
  - Pooling Layers
  - Nonlinearities

B. Leibe

## Recap: Recognition with Local Features

- Image content is transformed into local features that are invariant to translation, rotation, and scale
- Goal: Verify if they belong to a consistent configuration



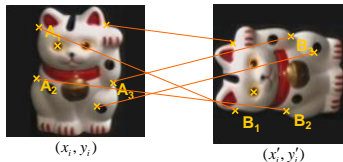
Local Features, e.g. SIFT

B. Leibe

Slide credit: David Lowe

## Recap: Fitting an Affine Transformation

- Assuming we know the correspondences, how do we get the transformation?



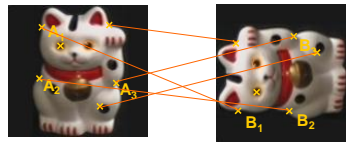
$$\begin{bmatrix} x'_i \\ y'_i \end{bmatrix} = \begin{bmatrix} m_1 & m_2 \\ m_3 & m_4 \end{bmatrix} \begin{bmatrix} x_i \\ y_i \end{bmatrix} + \begin{bmatrix} t_1 \\ t_2 \end{bmatrix}$$

$$\begin{bmatrix} x_i & y_i & 0 & 0 & 1 & 0 \\ 0 & 0 & x_i & y_i & 0 & 1 \\ & & & & & t_1 \\ & & & & & t_2 \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ m_3 \\ m_4 \\ t_1 \\ t_2 \end{bmatrix} = \begin{bmatrix} x'_i \\ y'_i \\ \dots \end{bmatrix}$$

B. Leibe

## Recap: Fitting a Homography

- Estimating the transformation



Homogenous coordinates:  $\begin{bmatrix} x' \\ y' \\ z' \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$

Image coordinates:  $\begin{bmatrix} x'' \\ y'' \\ z'' \end{bmatrix} = \begin{bmatrix} x'' \\ y'' \\ 1 \end{bmatrix}$

Matrix notation:  $x' = Hx$ ,  $x'' = \frac{1}{z'} x'$

$$x_A = \frac{h_{11}x_{B1} + h_{12}y_{B1} + h_{13}}{h_{31}x_{B1} + h_{32}y_{B1} + 1}$$

$$y_A = \frac{h_{21}x_{B1} + h_{22}y_{B1} + h_{23}}{h_{31}x_{B1} + h_{32}y_{B1} + 1}$$

B. Leibe

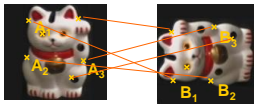
Slide credit: Krystian Mikolajczyk

Computer Vision Summer'19 RWTH AACHEN UNIVERSITY

### Recap: Fitting a Homography

- Estimating the transformation

$$\begin{aligned} h_{11}x_{A_1} + h_{12}y_{A_1} + h_{13} - x_{A_1}h_{31} - x_{A_1}h_{32}y_{A_1} - x_{A_1} &= 0 \\ h_{21}x_{A_1} + h_{22}y_{A_1} + h_{23} - y_{A_1}h_{31} - y_{A_1}h_{32}y_{A_1} - y_{A_1} &= 0 \end{aligned}$$



$$\begin{bmatrix} x_{A_1} & y_{A_1} & 1 & 0 & 0 & 0 & -x_{A_1}x_{B_1} & -x_{A_1}y_{B_1} & -x_{A_1} \\ 0 & 0 & 0 & x_{A_1} & y_{A_1} & 1 & -y_{A_1}x_{B_1} & -y_{A_1}y_{B_1} & -y_{A_1} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} h_{11} \\ h_{12} \\ h_{13} \\ h_{21} \\ h_{22} \\ h_{23} \\ h_{31} \\ h_{32} \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \end{bmatrix}$$

$Ah = 0$

Slide credit: Krystian Mikolajczyk B. Leibe 7

Computer Vision Summer'19 RWTH AACHEN UNIVERSITY

### Recap: Fitting a Homography

- Estimating the transformation
- Solution:
  - Null-space vector of A
  - Corresponds to smallest eigenvector

SVD

$$Ah = 0$$

$$A = UDV^T = U \begin{bmatrix} d_{11} & \dots & d_{19} \\ \vdots & \ddots & \vdots \\ d_{91} & \dots & d_{99} \end{bmatrix} \begin{bmatrix} v_{11} & \dots & v_{19} \\ \vdots & \ddots & \vdots \\ v_{91} & \dots & v_{99} \end{bmatrix}^T$$

$$h = [v_{19}, \dots, v_{99}]$$

Minimizes least square error

Slide credit: Krystian Mikolajczyk B. Leibe 8

Computer Vision Summer'19 RWTH AACHEN UNIVERSITY

### Recap: A General Point

- Equations of the form
 
$$Ax = 0$$

Think of this as an eigenvector equation  $Ax = \lambda x$  for the special case of  $\lambda = 0$ .
- How do we solve them? (always!)
  - SVD is the generalization of the eigenvector decomposition for non-square matrices A.
  - Apply SVD

$$A = UDV^T = U \begin{bmatrix} d_{11} & & & \\ & \ddots & & \\ & & d_{NN} & \\ & & & \vdots \\ & & & & v_{N1} & \dots & v_{NN} \end{bmatrix}^T$$

Singular values      Singular vectors

- Singular values of A = square roots of the eigenvalues of  $A^T A$ .
- The solution of  $Ax=0$  is the *nullspace* vector of A.
- This corresponds to the *smallest singular vector* of A.

B. Leibe 9

Computer Vision Summer'19 RWTH AACHEN UNIVERSITY

### Recap: Object Recognition by Alignment

- Assumption
  - Known object, rigid transformation compared to model image
  - $\Rightarrow$  If we can find evidence for such a transformation, we have recognized the object.
- You learned methods for
  - Fitting an *affine transformation* from  $\geq 3$  correspondences
  - Fitting a *homography* from  $\geq 4$  correspondences

Affine: solve a system

$$At = b$$

Homography: solve a system

$$Ah = 0$$

- Correspondences may be noisy and may contain outliers
  - $\Rightarrow$  Need to use robust methods that can filter out outliers

B. Leibe 10

Computer Vision Summer'19 RWTH AACHEN UNIVERSITY

### Topics of This Lecture

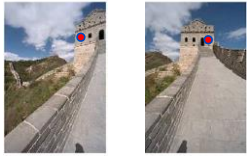

- Recap: Recognition with Local Features
- Dealing with Outliers
  - RANSAC
  - Generalized Hough Transform
- Deep Learning
  - Motivation
  - Neural Networks
- Convolutional Neural Networks
  - Convolutional Layers
  - Pooling Layers
  - Nonlinearities

B. Leibe 11

Computer Vision Summer'19 RWTH AACHEN UNIVERSITY

### Problem: Outliers

- Outliers can hurt the quality of our parameter estimates, e.g.,
  - An erroneous pair of matching points from two images
  - A feature point that is noise or doesn't belong to the transformation we are fitting.

Slide credit: Kristen Grauman B. Leibe 12

RWTH AACHEN UNIVERSITY

### Example: Least-Squares Line Fitting

- Assuming all the points that belong to a particular line are known

13  
Source: Forsyth & Ponce

RWTH AACHEN UNIVERSITY

### Outliers Affect Least-Squares Fit

14  
Source: Forsyth & Ponce

RWTH AACHEN UNIVERSITY

### Outliers Affect Least-Squares Fit

15  
Source: Forsyth & Ponce

RWTH AACHEN UNIVERSITY

### Strategy 1: RANSAC [Fischler81]

- RANdom SAmples Consensus
- Approach: we want to avoid the impact of outliers, so let's look for "inliers", and use only those.
- Intuition: if an outlier is chosen to compute the current fit, then the resulting line won't have much support from rest of the points.

16  
Slide credit: Kristen Grauman

RWTH AACHEN UNIVERSITY

### RANSAC

RANSAC loop:

- Randomly select a *seed group* of points on which to base transformation estimate (e.g., a group of matches)
- Compute transformation from seed group
- Find *inliers* to this transformation
- If the number of inliers is sufficiently large, re-compute least-squares estimate of transformation on all of the inliers

- Keep the transformation with the largest number of inliers

17  
Slide credit: Kristen Grauman

RWTH AACHEN UNIVERSITY

### RANSAC Line Fitting Example

- Task: Estimate the best line
  - How many points do we need to estimate the line?

18  
Slide credit: Jinxiang Chai

RWTH AACHEN UNIVERSITY

### RANSAC Line Fitting Example

- Task: Estimate the best line

Sample two points

19

Computer Vision Summer'19  
Slide credit: Jinxiang Chai  
B. Leibe

RWTH AACHEN UNIVERSITY

### RANSAC Line Fitting Example

- Task: Estimate the best line

Fit a line to them

20

Computer Vision Summer'19  
Slide credit: Jinxiang Chai  
B. Leibe

RWTH AACHEN UNIVERSITY

### RANSAC Line Fitting Example

- Task: Estimate the best line

Total number of points within a threshold of line.

21

Computer Vision Summer'19  
Slide credit: Jinxiang Chai  
B. Leibe

RWTH AACHEN UNIVERSITY

### RANSAC Line Fitting Example

- Task: Estimate the best line

"7 inlier points"  
Total number of points within a threshold of line.

22

Computer Vision Summer'19  
Slide credit: Jinxiang Chai  
B. Leibe

RWTH AACHEN UNIVERSITY

### RANSAC Line Fitting Example

- Task: Estimate the best line

Repeat, until we get a good result.

23

Computer Vision Summer'19  
Slide credit: Jinxiang Chai  
B. Leibe

RWTH AACHEN UNIVERSITY

### RANSAC Line Fitting Example

- Task: Estimate the best line

"11 inlier points"  
Repeat, until we get a good result.

24

Computer Vision Summer'19  
Slide credit: Jinxiang Chai  
B. Leibe

RWTH AACHEN UNIVERSITY

## RANSAC: How many samples?

- How many samples are needed?
  - Suppose  $w$  is fraction of inliers (points from line).
  - $n$  points needed to define hypothesis (2 for lines)
  - $k$  samples chosen.
- Prob. that a single sample of  $n$  points is correct:  $w^n$
- Prob. that all  $k$  samples fail is:  $(1 - w^n)^k$

⇒ Choose  $k$  high enough to keep this below the desired failure rate.

25

RWTH AACHEN UNIVERSITY

## RANSAC: Computed k (p=0.99)

Sample size n	Proportion of outliers						
	5%	10%	20%	25%	30%	40%	50%
2	2	3	5	6	7	11	17
3	3	4	7	9	11	19	35
4	3	5	9	13	17	34	72
5	4	6	12	17	26	57	146
6	4	7	16	24	37	97	293
7	4	8	20	33	54	163	588
8	5	9	26	44	78	272	1177

26

RWTH AACHEN UNIVERSITY

## After RANSAC

- RANSAC divides data into inliers and outliers and yields estimate computed from minimal set of inliers.
- Improve this initial estimate with estimation over all inliers (e.g. with standard least-squares minimization).
- But this may change inliers, so alternate fitting with re-classification as inlier/outlier.

27

RWTH AACHEN UNIVERSITY

## Example: Finding Feature Matches

- Find best stereo match within a square search window (here 300 pixels<sup>2</sup>)
- Global transformation model: epipolar geometry

28

RWTH AACHEN UNIVERSITY

## Example: Finding Feature Matches

- Find best stereo match within a square search window (here 300 pixels<sup>2</sup>)
- Global transformation model: epipolar geometry

before RANSAC

after RANSAC

29

RWTH AACHEN UNIVERSITY

## Problem with RANSAC

- In many practical situations, the percentage of outliers (incorrect putative matches) is often very high (90% or above).
- Alternative strategy: Generalized Hough Transform

30

RWTH AACHEN UNIVERSITY

## Strategy 2: Generalized Hough Transform

- Suppose our features are scale- and rotation-invariant
  - Then a single feature match provides an alignment hypothesis (translation, scale, orientation).

model

Computer Vision Summer'19 31

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

## Strategy 2: Generalized Hough Transform

- Suppose our features are scale- and rotation-invariant
  - Then a single feature match provides an alignment hypothesis (translation, scale, orientation).
  - Of course, a hypothesis from a single match is unreliable.
  - Solution: let each match vote for its hypothesis in a Hough space with very coarse bins.

model

Computer Vision Summer'19 32

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Recap: Recognition with Local Features
  - Dealing with Outliers
    - RANSAC
    - Generalized Hough Transform
- Deep Learning
  - Motivation
  - Neural Networks
- Convolutional Neural Networks
  - Convolutional Layers
  - Pooling Layers
  - Nonlinearities

Computer Vision Summer'19 33

B. Leibe

RWTH AACHEN UNIVERSITY

## We've finally got there!

# Deep Learning

Computer Vision Summer'19 34

B. Leibe

RWTH AACHEN UNIVERSITY

## Traditional Recognition Approach

Image/  
Video  
Pixels

Hand-designed  
feature  
extraction

Trainable  
classifier

Object  
Class

- Characteristics
  - Features are not learned, but engineered
  - Trainable classifier is often generic (e.g., SVM)
  - ⇒ Many successes in 2000-2010.

Computer Vision Summer'19 35

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

## Traditional Recognition Approach

- Features are key to recent progress in recognition
  - Multitude of hand-designed features currently in use
  - SIFT, HOG, .....
  - ⇒ Where next? Better classifiers? Or keep building more features?

DPM  
[Felzenszwalb  
et al., PAMI'07]

Dense SIFT+LBP+HOG → BOW → Classifier  
[Yan & Huan '10]  
(Winner of PASCAL 2010 Challenge)

Computer Vision Summer'19 36

Slide credit: Svetlana Lazebnik

RWTH AACHEN UNIVERSITY

## What About Learning the Features?

- Learn a *feature hierarchy* all the way from pixels to classifier
  - Each layer extracts features from the output of previous layer
  - Train all layers jointly

Image/ Video Pixels → Layer 1 → Layer 2 → Layer 3 → Simple Classifier

Computer Vision Summer'19 37

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

## “Shallow” vs. “Deep” Architectures

**Traditional recognition: “Shallow” architecture**

Image/ Video Pixels → Hand-designed feature extraction → Trainable classifier → Object Class

**Deep learning: “Deep” architecture**

Image/ Video Pixels → Layer 1 → ... → Layer N → Simple classifier → Object Class

Computer Vision Summer'19 38

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

## Background: Perceptrons

**Input**

**Weights**

$x_1$   $w_1$   
 $x_2$   $w_2$   
 $x_3$   $w_3$   
 $\vdots$   
 $x_d$   $w_d$

**Output:**  $\sigma(\mathbf{w} \cdot \mathbf{x} + \mathbf{b})$

**Sigmoid function**

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

Computer Vision Summer'19 39

Slide credit: Svetlana Lazebnik

RWTH AACHEN UNIVERSITY

## Inspiration: Neuron Cells

Computer Vision Summer'19 40

Slide credit: Svetlana Lazebnik, Rob Fergus

RWTH AACHEN UNIVERSITY

## Background: Multi-Layer Neural Networks

- Nonlinear classifier**
  - Training:** find network weights  $\mathbf{w}$  to minimize the error between true training labels  $t_n$  and estimated labels  $f_{\mathbf{w}}(x_n)$ :
 
$$E(\mathbf{W}) = \sum^n L(t_n, f(x_n; \mathbf{W}))$$
  - Minimization can be done by gradient descent, provided  $f$  is differentiable
    - Training method: **Error backpropagation.**

Computer Vision Summer'19 41

Slide credit: Svetlana Lazebnik B. Leibe

RWTH AACHEN UNIVERSITY

## Hubel/Wiesel Architecture

- D. Hubel, T. Wiesel (1959, 1962, Nobel Prize 1981)
  - Visual cortex consists of a hierarchy of *simple*, *complex*, and *hyper-complex* cells

**Hubel & Wiesel**  
topographical mapping

**featural hierarchy**

- hyper-complex cells
- complex cells
- simple cells

Legend:  $\square$  high level,  $\circ$  mid level,  $\bullet$  low level

Computer Vision Summer'19 42

Slide credit: Svetlana Lazebnik, Rob Fergus B. Leibe

**RWTH AACHEN UNIVERSITY**

## Convolutional Neural Networks (CNN, ConvNet)

- Neural network with specialized connectivity structure
  - Stack multiple stages of feature extractors
  - Higher stages compute more global, more invariant features
  - Classification layer at the end

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, [Gradient-based learning applied to document recognition](#), Proceedings of the IEEE 86(11): 2278–2324, 1998.

Slide credit: Svetlana Lazebnik

B. Leibe 43

**RWTH AACHEN UNIVERSITY**

## Topics of This Lecture

- Recap: Recognition with Local Features
  - RANSAC
  - Generalized Hough Transform
- Deep Learning
  - Motivation
  - Neural Networks
- Convolutional Neural Networks
  - Convolutional Layers
  - Pooling Layers
  - Nonlinearities

Computer Vision Summer'19

B. Leibe 45

**RWTH AACHEN UNIVERSITY**

## Convolutional Networks: Structure

- Feed-forward feature extraction
  - Conolve input with learned filters
  - Non-linearity
  - Spatial pooling
  - (Normalization)
- Supervised training of convolutional filters by back-propagating classification error

Computer Vision Summer'19

Slide credit: Svetlana Lazebnik

B. Leibe 46

**RWTH AACHEN UNIVERSITY**

## Convolutional Networks: Intuition

- Fully connected network
  - E.g. 1000×1000 image
  - 1M hidden units
  - ⇒ 1T parameters!
- Ideas to improve this
  - Spatial correlation is local

Computer Vision Summer'19

Slide adapted from Marc'Aurelio Ranzato

B. Leibe 47

**RWTH AACHEN UNIVERSITY**

## Convolutional Networks: Intuition

- Locally connected net
  - E.g. 1000×1000 image
  - 1M hidden units
  - 10×10 receptive fields
  - ⇒ 100M parameters!
- Ideas to improve this
  - Spatial correlation is local
  - Want translation invariance

Computer Vision Summer'19

Slide adapted from Marc'Aurelio Ranzato

B. Leibe 48

**RWTH AACHEN UNIVERSITY**

## Convolutional Networks: Intuition

- Convolutional net
  - Share the same parameters across different locations
  - Convolutions with learned kernels

Computer Vision Summer'19

Slide adapted from Marc'Aurelio Ranzato

B. Leibe 49



Computer Vision Summer'19

## Convolutional Networks: Intuition

RWTH AACHEN UNIVERSITY

- Convolutional net
  - Share the same parameters across different locations
  - Convolutions with learned kernels
- Learn *multiple* filters
  - E.g.  $1000 \times 1000$  image
  - 100 filters
  - $10 \times 10$  filter size
  - ⇒ 10k parameters
- Result: Response map
  - size:  $1000 \times 1000 \times 100$
  - Only memory, not params!

Slide adapted from Marc'Aurelio Ranzato. B. Leibe. Image source: Yann Lecun. 50

Computer Vision Summer'19

## Important Conceptual Shift

RWTH AACHEN UNIVERSITY

- Before
  - Diagram: A fully connected input layer (3 nodes) connects to a fully connected hidden layer (4 nodes), which connects to a fully connected output layer (2 nodes).
- Now:
  - Diagram: A 3D volume (red cube) is processed to produce another 3D volume (blue cube), which is then processed to produce a final 3D volume (green cube).

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe. 51

Computer Vision Summer'19

## Convolution Layers

RWTH AACHEN UNIVERSITY

Example image:  $32 \times 32 \times 3$  volume

Before: Full connectivity  $32 \times 32 \times 3$  weights

Now: Local connectivity  
One neuron connects to, e.g.,  $5 \times 5 \times 3$  region.  
⇒ Only  $5 \times 5 \times 3$  shared weights.

- Note: Connectivity is
  - Local in space ( $5 \times 5$  inside  $32 \times 32$ )
  - But full in depth (all 3 depth channels)

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe. 52

Computer Vision Summer'19

## Convolution Layers

RWTH AACHEN UNIVERSITY

before: "hidden layer of 200 neurons"  
now: "output volume of depth 200"

- All Neural Net activations arranged in 3 dimensions
  - Multiple neurons all looking at the same input region, stacked in depth

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe. 53

Computer Vision Summer'19

## Convolution Layers

RWTH AACHEN UNIVERSITY

Naming convention:

HEIGHT  
WIDTH  
DEPTH

- All Neural Net activations arranged in 3 dimensions
  - Multiple neurons all looking at the same input region, stacked in depth
  - Form a single  $[1 \times 1 \times \text{depth}]$  depth column in output volume.

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe. 54

Computer Vision Summer'19

## Convolution Layers


RWTH AACHEN UNIVERSITY


- All Neural Net activations arranged in 3 dimensions
  - Convolution layers can be stacked
  - The filters of the next layer then operate on the full activation volume.
  - Filters are local in (x,y), but densely connected in depth.

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe. 55

RWTH AACHEN UNIVERSITY

## Activation Maps of Convolutional Filters

Activations:  one filter = one depth slice (or activation map) 5 × 5 filters

Activations:  Each activation map is a depth slice through the output volume.

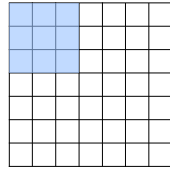
Activation maps

Computer Vision Summer'19 56

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

## Convolution Layers



Example:  
7 × 7 input  
assume 3 × 3 connectivity  
stride 1

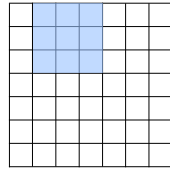
- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 57

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

## Convolution Layers



Example:  
7 × 7 input  
assume 3 × 3 connectivity  
stride 1

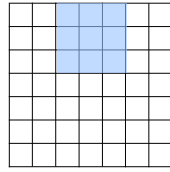
- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 58

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

## Convolution Layers



Example:  
7 × 7 input  
assume 3 × 3 connectivity  
stride 1

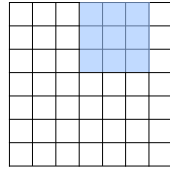
- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 59

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

## Convolution Layers



Example:  
7 × 7 input  
assume 3 × 3 connectivity  
stride 1

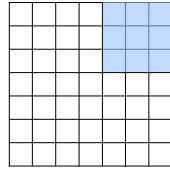
- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 60

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

## Convolution Layers



Example:  
7 × 7 input  
assume 3 × 3 connectivity  
stride 1  
⇒ 5 × 5 output

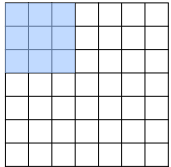
- Replicate this column of hidden neurons across space, with some **stride**.

Computer Vision Summer'19 61

Slide credit: FeiFei Li, Andrei Karpathy. B. Leibe

RWTH AACHEN UNIVERSITY

## Convolution Layers



Example:  
 $7 \times 7$  input  
 assume  $3 \times 3$  connectivity  
 stride 1  
 $\Rightarrow 5 \times 5$  output

What about stride 2?

- Replicate this column of hidden neurons across space, with some **stride**.

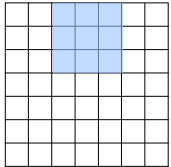
Computer Vision Summer'19

62

Slide credit: FeiFei Li, Andrej Karpathy B. Leibe

RWTH AACHEN UNIVERSITY

## Convolution Layers



Example:  
 $7 \times 7$  input  
 assume  $3 \times 3$  connectivity  
 stride 1  
 $\Rightarrow 5 \times 5$  output

What about stride 2?

- Replicate this column of hidden neurons across space, with some **stride**.

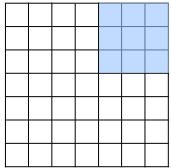
Computer Vision Summer'19

63

Slide credit: FeiFei Li, Andrej Karpathy B. Leibe

RWTH AACHEN UNIVERSITY

## Convolution Layers



Example:  
 $7 \times 7$  input  
 assume  $3 \times 3$  connectivity  
 stride 1  
 $\Rightarrow 5 \times 5$  output

What about stride 2?  
 $\Rightarrow 3 \times 3$  output

- Replicate this column of hidden neurons across space, with some **stride**.

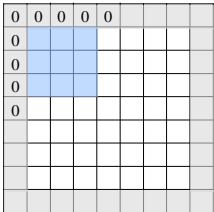
Computer Vision Summer'19

64

Slide credit: FeiFei Li, Andrej Karpathy B. Leibe

RWTH AACHEN UNIVERSITY

## Convolution Layers



Example:  
 $7 \times 7$  input  
 assume  $3 \times 3$  connectivity  
 stride 1  
 $\Rightarrow 5 \times 5$  output

What about stride 2?  
 $\Rightarrow 3 \times 3$  output

- Replicate this column of hidden neurons across space, with some **stride**.
- In practice, common to zero-pad the border.
  - Preserves the size of the input spatially.

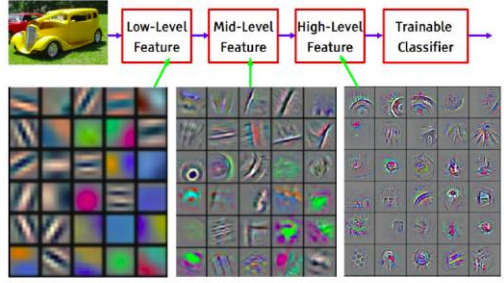
Computer Vision Summer'19

65

Slide credit: FeiFei Li, Andrej Karpathy B. Leibe

RWTH AACHEN UNIVERSITY

## Effect of Multiple Convolution Layers



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

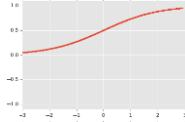
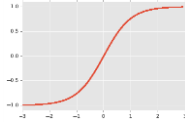
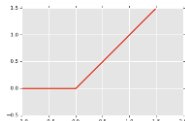
Computer Vision Summer'19

66

Slide credit: Yann LeCun B. Leibe

RWTH AACHEN UNIVERSITY

## Commonly Used Nonlinearities

- Sigmoid
 
$$g(a) = \sigma(a) = \frac{1}{1 + \exp\{-a\}}$$

- Hyperbolic tangent
 
$$g(a) = \tanh(a) = 2\sigma(2a) - 1$$

- Rectified linear unit (ReLU)
 
$$g(a) = \max\{0, a\}$$


Preferred option for deep networks

Computer Vision Summer'19

67

B. Leibe

Computer Vision Summer'19

## Convolutional Networks: Intuition

RWTH AACHEN UNIVERSITY

- Let's assume the filter is an eye detector
  - How can we make the detection robust to the exact location of the eye?

Slide adapted from Marc'Aurelio Ranzato. B. Leibe. Image source: Yann Lecun. 68

Computer Vision Summer'19

## Convolutional Networks: Intuition

RWTH AACHEN UNIVERSITY

- Let's assume the filter is an eye detector
  - How can we make the detection robust to the exact location of the eye?
- Solution:
  - By **pooling** (e.g., max or avg) filter responses at different spatial locations, we gain robustness to the exact spatial location of features.

Slide adapted from Marc'Aurelio Ranzato. B. Leibe. Image source: Yann Lecun. 69

Computer Vision Summer'19

## Max Pooling

RWTH AACHEN UNIVERSITY

Single depth slice

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters and stride 2

6	8
3	4

- Effect:
  - Make the representation smaller without losing too much information
  - Achieve robustness to translations

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe. 70

Computer Vision Summer'19

## Max Pooling

RWTH AACHEN UNIVERSITY

Single depth slice

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters and stride 2

6	8
3	4

- Note
  - Pooling happens independently across each slice, preserving the number of slices.

Slide adapted from FeiFei Li, Andrei Karpathy. B. Leibe. 71

Computer Vision Summer'19

## Compare: SIFT Descriptor

RWTH AACHEN UNIVERSITY

Image Pixels

Apply oriented filters [Lowe [ICV 2004]]

Spatial pool (Sum)

Normalize to unit length

Feature Vector

Slide credit: Svetlana Lazebnik. B. Leibe. 72

Computer Vision Summer'19

## Compare: Spatial Pyramid Matching

RWTH AACHEN UNIVERSITY

SIFT features

Filter with Visual Words [Lazebnik, Schmid, Ponce [CVPR 2006]]

Take max VW response (L-inf normalization)

Multi-scale spatial pool (Sum)

Global image descriptor

Slide credit: Svetlana Lazebnik. B. Leibe. 73

## References and Further Reading

- More information on Deep Learning and CNNs can be found in Chapters 6 and 9 of the Goodfellow & Bengio book

I. Goodfellow, Y. Bengio, A. Courville  
Deep Learning  
MIT Press, 2016  
<http://www.deeplearningbook.org/>

