

RWTH AACHEN
UNIVERSITY

Advanced Machine Learning Lecture 16

Convolutional Neural Networks II

14.01.2016

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

Advanced Machine Learning Winter'15

RWTH AACHEN
UNIVERSITY

Announcements

- Lecture evaluation
 - Please fill out the evaluation forms.

B. Leibe

2

Advanced Machine Learning Winter'15

RWTH AACHEN
UNIVERSITY

This Lecture: *Advanced Machine Learning*

- Regression Approaches
 - Linear Regression
 - Regularization (Ridge, Lasso)
 - Gaussian Processes
- Learning with Latent Variables
 - Prob. Distributions & Approx. Inference
 - Mixture Models
 - EM and Generalizations
- Deep Learning
 - Linear Discriminants
 - Neural Networks
 - Backpropagation & Optimization
 - CNNs, RNNs, RBMs, etc.

B. Leibe

Advanced Machine Learning Winter'15

RWTH AACHEN
UNIVERSITY

Topics of This Lecture

- Recap: CNNs
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
- Visualizing CNNs
 - Visualizing CNN features
 - Visualizing responses
 - Visualizing learned structures
- Applications

B. Leibe

4

Advanced Machine Learning Winter'15

RWTH AACHEN
UNIVERSITY

Recap: Convolutional Neural Networks

- 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.

B. Leibe

5

Advanced Machine Learning Winter'15

RWTH AACHEN
UNIVERSITY

Recap: Intuition of CNNs

- Convolutional net
 - Share the same parameters across different locations
 - Convolutions with learned kernels
- Learn *multiple* filters
 - E.g. 1000x1000 image
 - 100 filters
 - 10x10 filter size
 - ⇒ only 10k parameters
- Result: Response map
 - size: 1000x1000x100
 - Only memory, not params!

B. Leibe

6

Advanced Machine Learning Winter'15

Advanced Machine Learning Winter'15

Recap: Convolution Layers

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

Advanced Machine Learning Winter'15

Recap: Activation Maps

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

5x5 filters

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

Activation maps

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

Advanced Machine Learning Winter'15

Recap: Pooling Layers

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

Advanced Machine Learning Winter'15

Topics of This Lecture

- Recap: CNNs
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
- Visualizing CNNs
 - Visualizing CNN features
 - Visualizing responses
 - Visualizing learned structures
- Applications

B. Leibe

Advanced Machine Learning Winter'15

Recap: ImageNet Challenge 2012

ImageNet

- ~14M labeled internet images
- 20k classes
- Human labels via Amazon Mechanical Turk

Challenge (ILSVRC)

- 1.2 million training images
- 1000 classes
- Goal: Predict ground-truth class within top-5 responses
- Currently one of the top benchmarks in Computer Vision

[Deng et al., CVPR'09]

B. Leibe

Advanced Machine Learning Winter'15

CNN Architectures: AlexNet (2012)

- Similar framework as LeNet, but
 - Bigger model (7 hidden layers, 650k units, 60M parameters)
 - More data (10^6 images instead of 10^3)
 - GPU implementation
 - Better regularization and up-to-date tricks for training (Dropout)

A. Krizhevsky, I. Sutskever, and G. Hinton, [ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012.

Image source: A. Krizhevsky, I. Sutskever and G.F. Hinton, NIPS 2012

Advanced Machine Learning Winter'15

ILSVRC 2012 Results

Top-5 error rate %

- AlexNet almost halved the error rate
 - 16.4% error (top-5) vs. 26.2% for the next best approach
 - ⇒ A revolution in Computer Vision
 - Acquired by Google in Jan '13, deployed in Google+ in May '13

B. Leibe 14

Advanced Machine Learning Winter'15

CNN Architectures: VGGNet (2015)

- Main ideas
 - Deeper network
 - Stacked convolutional layers with smaller filters (+ nonlinearity)
 - Detailed evaluation of all components

ConvNet Configurations				
A	A-LRN	B	C	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	19 weight layers
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
		input (224 × 224 RGB image)		
		maxpool		
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
		maxpool		
conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256
		maxpool		
conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512
		maxpool		
conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512
		maxpool		
		maxpool		
		FC-4096		
		FC-4096		
		FC-1000		
		softmax		

B. Leibe 16

Advanced Machine Learning Winter'15

Comparison: AlexNet vs. VGGNet

Legend:

- Input: Image input
- Conv: Convolutional layer
- Pool: Max-pooling layer
- FC: Fully-connected layer
- Softmax: Softmax layer

K. Simonyan, A. Zisserman, [Very Deep Convolutional Networks for Large-Scale Image Recognition](#), ICLR 2015

B. Leibe 17

Advanced Machine Learning Winter'15

Comparison: AlexNet vs. VGGNet

- Receptive fields in the first layer
 - AlexNet: 11 × 11, stride 4
 - Zeiler & Fergus: 7 × 7, stride 2
 - VGGNet: 3 × 3, stride 1
- Why that?
 - If you stack three 3 × 3 on top of another 3 × 3 layer, you effectively get a 5 × 5 receptive field.
 - With three 3 × 3 layers, the receptive field is already 7 × 7.
 - But much fewer parameters: $3 \cdot 3^2 = 27$ instead of $7^2 = 49$.
 - In addition, non-linearities in-between 3 × 3 layers for additional discriminativity.

B. Leibe 18

Advanced Machine Learning Winter'15

CNN Architectures: GoogLeNet (2014)

Inception module with dimension reductions

- Main ideas
 - “Inception” module as modular component
 - Learns filters at several scales within each module

C. Szegedy, W. Liu, Y. Jia, et al, [Going Deeper with Convolutions](#), arXiv:1409.4842, 2014.

B. Leibe 20

Advanced Machine Learning Winter'15

GoogLeNet Visualization

Inception module + copies

Auxiliary classification outputs for training the lower layers (deprecated)

Convolution
Pooling
Softmax
Other

B. Leibe 21

Results on ILSVRC

Method	top-1 val. error (%)	top-5 val. error (%)	top-5 test error (%)
VGG (2 nets, multi-crop & dense eval.)	23.7	6.8	6.8
VGG (1 net, multi-crop & dense eval.)	24.4	7.1	7.0
VGG (ILSVRC submission, 7 nets, dense eval.)	24.7	7.5	7.3
GoogLeNet (Szegedy et al., 2014) (1 net)	-	-	7.9
GoogLeNet (Szegedy et al., 2014) (7 nets)	-	-	6.7
MSRA (He et al., 2014) (11 nets)	-	-	8.1
MSRA (He et al., 2014) (1 net)	27.9	9.1	9.1
Clarifai (Russakovsky et al., 2014) (multiple nets)	-	-	11.7
Clarifai (Russakovsky et al., 2014) (1 net)	-	-	12.5
Zeiler & Fergus (Zeiler & Fergus, 2013) (6 nets)	36.0	14.7	14.8
Zeiler & Fergus (Zeiler & Fergus, 2013) (1 net)	37.5	16.0	16.1
OverFeat (Sermanet et al., 2014) (7 nets)	34.0	13.2	13.6
OverFeat (Sermanet et al., 2014) (1 net)	35.7	14.2	-
Krizhevsky et al. (Krizhevsky et al., 2012) (5 nets)	38.1	16.4	16.4
Krizhevsky et al. (Krizhevsky et al., 2012) (1 net)	40.7	18.2	-


- VGGNet and GoogLeNet perform at similar level
 - Comparison: human performance ~5% [Karpathy]

<http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-cornet-on-imagenet/>

B. Leibe 22

Understanding the ILSVRC Challenge

- Imagine the scope of the problem!
 - 1000 categories
 - 1.2M training images
 - 50k validation images
- This means...
 - Speaking out the list of category names at 1 word/s...
 - ...takes 15mins.
 - Watching a slideshow of the validation images at 2s/image...
 - ...takes a full day (24h+).
 - Watching a slideshow of the training images at 2s/image...
 - ...takes a full month.




B. Leibe 23

More Finegrained Classes



B. Leibe 25

Quirks and Limitations of the Data Set



- Generated from WordNet ontology
 - Some animal categories are overrepresented
 - E.g., 120 subcategories of dog breeds

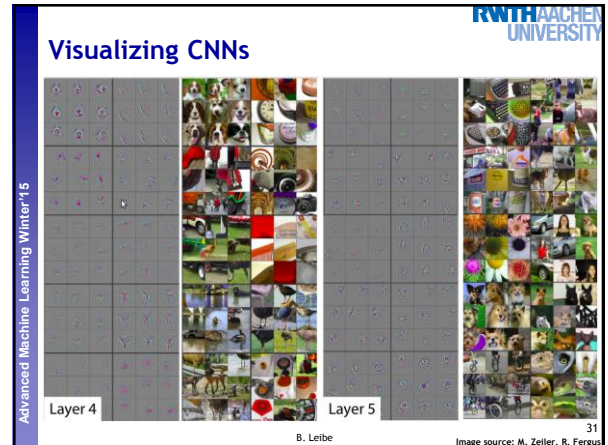
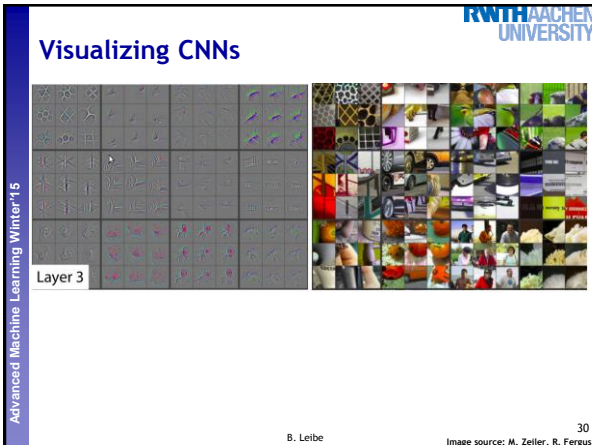
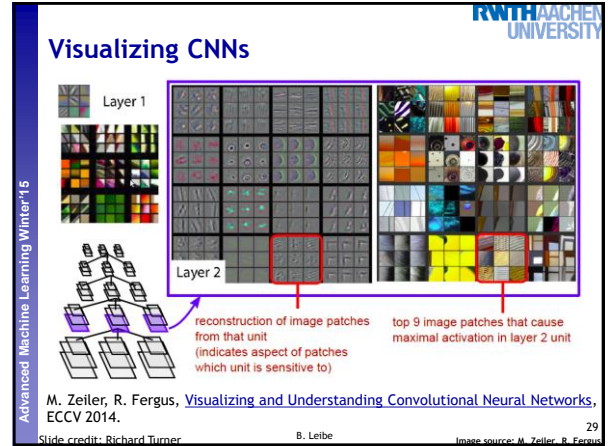
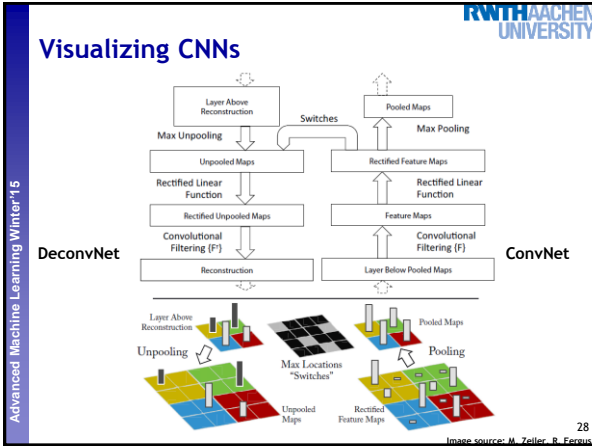
⇒ 6.7% top-5 error looks all the more impressive

B. Leibe 26

Topics of This Lecture

- Recap: CNNs
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
- Visualizing CNNs
 - Visualizing CNN features
 - Visualizing responses
 - Visualizing learned structures
- Applications

B. Leibe 27



Advanced Machine Learning Winter'15

What Does the Network React To?

- Occlusion Experiment
 - Mask part of the image with an occluding square.
 - Monitor the output

B. Leibe

32

Advanced Machine Learning Winter'15

What Does the Network React To?

Input image

True Label: Pomeranian

p(True class)

Most probable class

Legend: Pomeranian, Tennis ball, Keeshond, Pekingese

Slide credit: Svetlana Lazebnik, Rob Fergus


Image source: M. Zeiler, R. Fergus

33

RWTH AACHEN UNIVERSITY

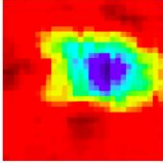
What Does the Network React To?

Input image

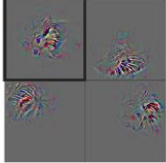


True Label: Pomeranian

Total activation in most active 5th layer feature map



Other activations from the same feature map.




Advanced Machine Learning Winter'15 34

Slide credit: Svetlana Lazebnik, Rob Fergus Image source: M. Zeiler, R. Fergus

RWTH AACHEN UNIVERSITY

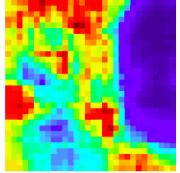
What Does the Network React To?

Input image

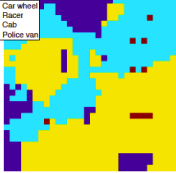


True Label: Car Wheel

p(True class)



Most probable class




Advanced Machine Learning Winter'15 35

Slide credit: Svetlana Lazebnik, Rob Fergus Image source: M. Zeiler, R. Fergus

RWTH AACHEN UNIVERSITY

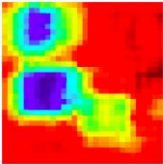
What Does the Network React To?

Input image




True Label: Car Wheel

Total activation in most active 5th layer feature map



Other activations from the same feature map.




Advanced Machine Learning Winter'15 36

Slide credit: Svetlana Lazebnik, Rob Fergus Image source: M. Zeiler, R. Fergus

RWTH AACHEN UNIVERSITY

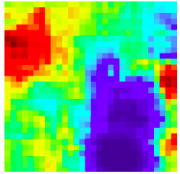
What Does the Network React To?

Input image

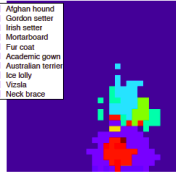


True Label: Afghan Hound

p(True class)



Most probable class




Advanced Machine Learning Winter'15 37

Slide credit: Svetlana Lazebnik, Rob Fergus Image source: M. Zeiler, R. Fergus

RWTH AACHEN UNIVERSITY

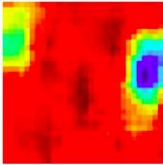
What Does the Network React To?

Input image




True Label: Afghan Hound

Total activation in most active 5th layer feature map



Other activations from the same feature map.




Advanced Machine Learning Winter'15 38


Slide credit: Svetlana Lazebnik, Rob Fergus Image source: M. Zeiler, R. Fergus

RWTH AACHEN UNIVERSITY

Inceptionism: Dreaming ConvNets


→

optimize with prior

→


- Idea
 - Start with a random noise image.
 - Enhance the input image such as to enforce a particular response (e.g., banana).
 - Combine with prior constraint that image should have similar statistics as natural images.
- ⇒ Network hallucinates characteristics of the learned class.

http://googleresearch.blogspot.de/2015/06/inceptionism-going-deeper-into-neural.html

Advanced Machine Learning Winter'15 39

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

Inceptionism: Dreaming ConvNets

- Results



<http://googleresearch.blogspot.de/2015/07/deepdream-code-example-for-visualizing.html>

40

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

Inceptionism: Dreaming ConvNets



<https://www.youtube.com/watch?v=IREsx-xWQ0g>

41

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Recap: CNNs
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet
- Visualizing CNNs
 - Visualizing CNN features
 - Visualizing responses
 - Visualizing learned structures
- Applications

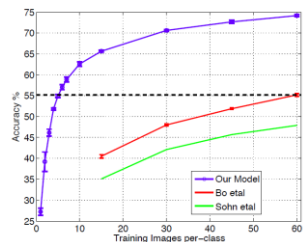
B. Leibe

42

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

The Learned Features are Generic



- Experiment: feature transfer
 - Train network on ImageNet
 - Chop off last layer and train classification layer on CalTech256

⇒ State of the art accuracy already with only 6 training images

B. Leibe

Image sources: M. Zeller, B. Feuz

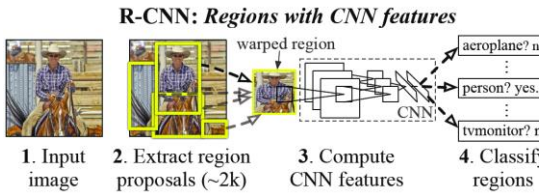
43

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

Other Tasks: Detection

R-CNN: Regions with CNN features



- Input image
- Extract region proposals (~2k)
- Compute CNN features
- Classify regions

Results on PASCAL VOC Detection benchmark

- Pre-CNN state of the art: 35.1% mAP [Uijlings et al., 2013]
- 33.4% mAP DPM
- R-CNN: 53.7% mAP

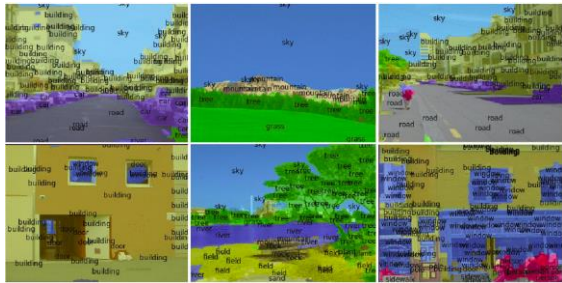
R. Girshick, J. Donahue, T. Darrell, and J. Malik, [Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation](#), CVPR 2014

44

Advanced Machine Learning Winter'15

RWTH AACHEN UNIVERSITY

Other Tasks: Semantic Segmentation



[Farabet et al. ICML 2012, PAMI 2013]

B. Leibe

45

RWTH AACHEN UNIVERSITY

Other Tasks: Semantic Segmentation

[Farabet et al. ICML 2012, PAMI 2013]

B. Leibe 46

RWTH AACHEN UNIVERSITY

Other Tasks: Face Verification

Y. Taigman, M. Yang, M. Ranzato, L. Wolf, **DeepFace: Closing the Gap to Human-Level Performance in Face Verification**, CVPR 2014

Slide credit: Svetlana Lazebnik 47

RWTH AACHEN UNIVERSITY

Commercial Recognition Services

- E.g., **clarifai**

- Be careful when taking test images from Google Search
 - Chances are they may have been seen in the training set...

B. Leibe 48

RWTH AACHEN UNIVERSITY

Commercial Recognition Services

B. Leibe 49

RWTH AACHEN UNIVERSITY

References and Further Reading

- LeNet**
 - Y. LeCun, L. Bottou, Y. Bengio, and P. Häffner, [Gradient-based learning applied to document recognition](#), Proceedings of the IEEE 86(11): 2278-2324, 1998.
- AlexNet**
 - A. Krizhevsky, I. Sutskever, and G. Hinton, [ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012.
- VGGNet**
 - K. Simonyan, A. Zisserman, [Very Deep Convolutional Networks for Large-Scale Image Recognition](#), ICLR 2015
- GoogLeNet**
 - C. Szegedy, W. Liu, Y. Jia, et al, [Going Deeper with Convolutions](#), arXiv:1409.4842, 2014.

B. Leibe 50

RWTH AACHEN UNIVERSITY

Effect of Multiple Convolution Layers

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

Slide credit: Yann LeCun 54