

RWTH AACHEN
UNIVERSITY

Machine Learning – Lecture 17

Convolutional Neural Networks III

08.01.2018

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

Machine Learning Winter '17

RWTH AACHEN
UNIVERSITY

Course Outline

- Fundamentals
 - Bayes Decision Theory
 - Probability Density Estimation
- Classification Approaches
 - Linear Discriminants
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Random Forests
- Deep Learning
 - Foundations
 - Convolutional Neural Networks**
 - Recurrent Neural Networks

B. Leibe

RWTH AACHEN
UNIVERSITY

Topics of This Lecture

- Recap: CNN Architectures
- Residual Networks
 - Detailed analysis
 - ResNets as ensembles of shallow networks
- Applications of CNNs
 - Object detection
 - Semantic segmentation
 - Face identification

B. Leibe

Machine Learning Winter '17

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

Machine Learning Winter '17

RWTH AACHEN
UNIVERSITY

Recap: 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.

B. Leibe

Machine Learning Winter '17

RWTH AACHEN
UNIVERSITY

Recap: VGGNet (2014/15)

- Main ideas
 - Deeper network
 - Stacked convolutional layers with smaller filters (+ nonlinearity)
 - Detailed evaluation of all components
- Results
 - Improved ILSVRC top-5 error rate to 6.7%.

| ConvNet Configuration | | | | | |
|-----------------------|------------------|--------------------------|--------------------------|--------------------------|--------------------------|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| conv-3-64 | conv-3-64 LRN | conv-3-64 conv-3-64 | conv-3-64 conv-3-64 | conv-3-64 conv-3-64 | conv-3-64 conv-3-64 |
| conv-3-128 | conv-3-128 | conv-3-128 conv-3-128 | conv-3-128 conv-3-128 | conv-3-128 conv-3-128 | conv-3-128 conv-3-128 |
| conv-3-256 | conv-3-256 | conv-3-256 | conv-3-256 conv-1-256 | conv-3-256 conv-3-256 | conv-3-256 conv-3-256 |
| conv-3-512 | conv-3-512 | conv-3-512 | conv-3-512 conv-1-512 | conv-3-512 conv-3-512 | conv-3-512 conv-3-512 |
| conv-3-512 | conv-3-512 | conv-3-512 | conv-3-512 conv-3-512 | conv-3-512 conv-3-512 | conv-3-512 conv-3-512 |
| conv-3-512 | conv-3-512 | conv-3-512 | conv-3-512 conv-3-512 | conv-3-512 conv-3-512 | conv-3-512 conv-3-512 |
| | | | conv-1-1000 | conv-1-1000 | conv-1-1000 |
| | | | FC-1000 | FC-1000 | FC-1000 |
| | | | self-max | self-max | self-max |
| | | | | Mainly used | |

B. Leibe

Machine Learning Winter '17

Machine Learning Winter '17

Recap: GoogLeNet (2014)

RWTH AACHEN UNIVERSITY

- Ideas:
 - Learn features at multiple scales
 - Modular structure

Inception module + copies

Auxiliary classification outputs for training the lower layers (deprecated)

Convolution
Pooling
Softmax
Other

(b) Inception module with dimension reductions

B. Leibe

Image source: Szepesvári et al.

8

Machine Learning Winter '17

Recap: Visualizing CNNs

RWTH AACHEN UNIVERSITY

Low-Level Feature → Mid-Level Feature → High-Level Feature → Trainable Classifier

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

Slide credit: Yann LeCun

B. Leibe

10

Machine Learning Winter '17

Topics of This Lecture

RWTH AACHEN UNIVERSITY

- Recap: CNN Architectures
- Residual Networks**
 - Detailed analysis
 - ResNets as ensembles of shallow networks
- Applications of CNNs
 - Object detection
 - Semantic segmentation
 - Face identification

B. Leibe

11

Machine Learning Winter '17

Recap: Residual Networks

RWTH AACHEN UNIVERSITY

AlexNet, 8 layers (ILSVRC 2012)

VGG, 19 layers (ILSVRC 2014)

GoogleNet, 22 layers (ILSVRC 2014)

Slide credit: Kaiming He

B. Leibe

12

Machine Learning Winter '17

Recap: Residual Networks

RWTH AACHEN UNIVERSITY

AlexNet, 8 layers (ILSVRC 2012)

VGG, 19 layers (ILSVRC 2014)

ResNet, 152 layers (ILSVRC 2015)

- Core component
 - Skip connections bypassing each layer
 - Better propagation of gradients to the deeper layers

$$F(x) = \text{weight layer} \rightarrow \text{relu}$$

$$H(x) = F(x) + x \rightarrow \text{relu}$$

B. Leibe

13

Machine Learning Winter '17

Spectrum of Depth

RWTH AACHEN UNIVERSITY

5 layers: easy

>10 layers: initialization, Batch Normalization

>30 layers: skip connections

>100 layers: identity skip connections

>1000 layers: ?

shallower ← → deeper

Slide credit: Kaiming He

B. Leibe

14

Spectrum of Depth

- Deeper models are more powerful
 - But training them is harder.
 - Main problem: getting the gradients back to the early layers
 - The deeper the network, the more effort is required for this.

Machine Learning Winter '17 | Slide adapted from Kaiming He | B. Leibe | 15

Initialization

22-layer ReLU net:
good init converges faster

30-layer ReLU net:
good init is able to converge

- Importance of proper initialization (Recall Lecture 14)
 - Glorot initialization for tanh nonlinearities
 - He initialization for ReLU nonlinearities
 - ⇒ For deep networks, this really makes a difference!

Machine Learning Winter '17 | Slide credit: Kaiming He | B. Leibe | 16

Batch Normalization

- Effect of batch normalization
 - Greatly improved speed of convergence

Machine Learning Winter '17 | Image source: Ioffe & Szegedy | B. Leibe | 17

Going Deeper

- Checklist
 - Initialization ok
 - Batch normalization ok
 - Are we now set?
 - Is learning better networks now as simple as stacking more layers?

Machine Learning Winter '17 | Slide credit: Kaiming He | B. Leibe | 18

Simply Stacking Layers?

CIFAR-10
train error (%)

CIFAR-10
test error (%)

- Experiment going deeper
 - Plain nets: stacking 3x3 convolution layers
 - ⇒ 56-layer net has higher training error than 20-layer net

Machine Learning Winter '17 | Slide credit: Kaiming He | B. Leibe | 19

Simply Stacking Layers?

CIFAR-10

ImageNet-1000

- General observation
 - Overly deep networks have higher training error
 - A general phenomenon, observed in many training sets

Machine Learning Winter '17 | Slide credit: Kaiming He | B. Leibe | 20

Why Is That???

- A deeper model should not have higher training error!
 - Richer solution space should allow it to find better solutions
- Solution by construction
 - Copy the original layers from a learned shallower model
 - Set the extra layers as identity
 - Such a network should achieve at least the same low training error.
- Reason: Optimization difficulties
 - Solvers cannot find the solution when going deeper...

Machine Learning Winter '17

Slide credit: Kaiqing He

B. Leibe

41

Deep Residual Learning

- Plain net
 - any two stacked layers
 - weight layer
 - relu
 - weight layer
 - relu
 - $H(x)$
- $H(x)$ is any desired mapping
- Hope the 2 weight layers fit $H(x)$

Machine Learning Winter '17

Slide credit: Kaiqing He

B. Leibe

22

Deep Residual Learning

- Residual net
 - $F(x)$
 - weight layer
 - relu
 - weight layer
 - identity x
 - $H(x) = F(x) + x$
 - relu
- $H(x)$ is any desired mapping
- Hope the 2 weight layers fit $H(x)$
- Hope the 2 weight layers fit $F(x)$
- Let $H(x) = F(x) + x$

Machine Learning Winter '17

Slide credit: Kaiqing He

B. Leibe

23

Deep Residual Learning

- $F(x)$ is a residual mapping w.r.t. identity
 - If identity were optimal, it is easy to set weights as 0
 - If optimal mapping is closer to identity, it is easier to find small fluctuations
 - Further advantage: direct path for the gradient to flow to the previous stages

Machine Learning Winter '17

Slide credit: Kaiqing He

B. Leibe

24

Network Design

- Simple, VGG-style design
 - (Almost) all 3x3 convolutions
 - Spatial size / 2 \Rightarrow #filters \cdot 2 (same complexity per layer)
 - Batch normalization
 - \Rightarrow Simple design, just deep.

Machine Learning Winter '17

B. Leibe

25

ImageNet Performance

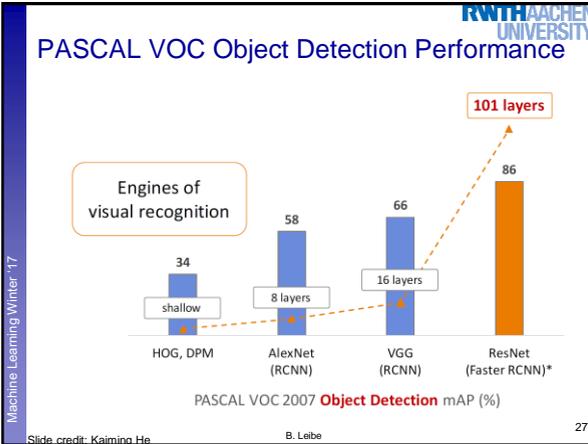
| Model | Layers | ImageNet Classification top-5 error (%) |
|---------------------|---------|---|
| ILSVRC'15 ResNet | 152 | 3.57 |
| ILSVRC'14 GoogleNet | 22 | 6.7 |
| ILSVRC'14 VGG | 19 | 7.3 |
| ILSVRC'13 AlexNet | 8 | 11.7 |
| ILSVRC'12 AlexNet | 8 | 16.4 |
| ILSVRC'11 shallow | shallow | 25.8 |
| ILSVRC'10 | shallow | 28.2 |

Machine Learning Winter '17

Slide credit: Kaiqing He

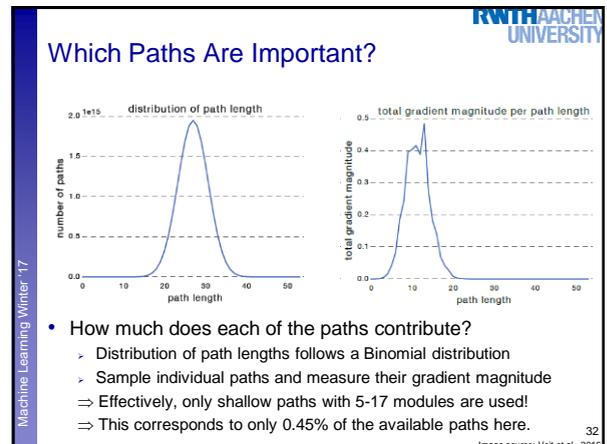
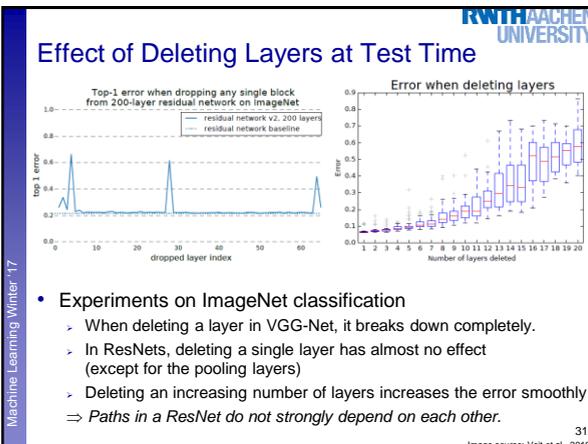
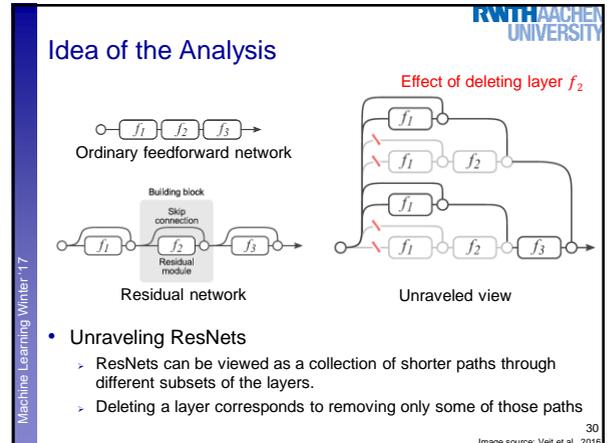
B. Leibe

26



- ### Topics of This Lecture
- Recap: CNN Architectures
 - Residual Networks
 - Detailed analysis
 - ResNets as ensembles of shallow networks
 - Applications of CNNs
 - Object detection
 - Semantic segmentation
 - Face identification

- ### What Is The Secret Behind ResNets?
- Empirically, they perform very well, but why is that?
 - He's original explanation [He, 2016]
 - ResNets allow gradients to pass through the skip connections in unchanged form.
 - This makes it possible to effectively train deeper networks.
 - ⇒ Secret of success: **depth is good**
 - More recent explanation [Veit, 2016]
 - ResNets actually do not use deep network paths.
 - Instead, they effectively implement an ensemble of shallow network paths.
 - ⇒ Secret of success: **ensembles are good**
- A. Veit, M. Wilber, S. Belongie, *Residual Networks Behave Like Ensembles of Relatively Shallow Networks*, NIPS 2016



Machine Learning Winter '17

Summary

- The effective paths in ResNets are relatively shallow
 - Effectively only 5-17 active modules
- This explains the resilience to deletion
 - Deleting any single layer only affects a subset of paths (and the shorter ones less than the longer ones).
- New interpretation of ResNets
 - ResNets work by creating an ensemble of relatively shallow paths
 - Making ResNets deeper increases the size of this ensemble
 - Excluding longer paths from training does not negatively affect the results.

Building block: skip connection, residual module, f_1 , f_2 , f_3

total gradient magnitude per path length

path length

33
Image source: Vait et al., 2014

Machine Learning Winter '17

Topics of This Lecture

- Recap: CNN Architectures
- Residual Networks
 - Detailed analysis
 - ResNets as ensembles of shallow networks
- Applications of CNNs
 - Object detection
 - Semantic segmentation
 - Face identification

B. Leibe

34

Machine Learning Winter '17

The Learned Features are Generic

Accuracy %

Training Images per-class

state of the art level (pre-CNN)

- Experiment: feature transfer
 - Train AlexNet-like network on ImageNet
 - Chop off last layer and train classification layer on CalTech256
 - State of the art accuracy already with only 6 training images!

35
B. Leibe
Image source: M. Zeller, R. Fergus

Machine Learning Winter '17

Transfer Learning with CNNs

1. Train on ImageNet
2. If small dataset: fix all weights (treat CNN as fixed feature extractor), retrain only the classifier

I.e., swap the Softmax layer at the end

36
Slide credit: Andrei Karpathy
B. Leibe

Machine Learning Winter '17

Transfer Learning with CNNs

1. Train on ImageNet
3. If you have medium sized dataset, "finetune" instead: use the old weights as initialization, train the full network or only some of the higher layers.

Retrain bigger portion of the network

37
B. Leibe

Machine Learning Winter '17

Other Tasks: Detection

R-CNN: Regions with CNN features

1. Input image
2. Extract region proposals (~2k)
3. Compute CNN features
4. 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

40

Machine Learning Winter '17

More Recent Version: Faster R-CNN

- One network, four losses
 - Remove dependence on external region proposal algorithm.
 - Instead, infer region proposals from same CNN.
 - Feature sharing
 - Joint training
 - Object detection in a single pass becomes possible.

Slide credit: Ross Girshick

41

Machine Learning Winter '17

Faster R-CNN (based on ResNets)

K. He, X. Zhang, S. Ren, J. Sun, [Deep Residual Learning for Image Recognition](#), CVPR 2016.

B. Leibe

42

Machine Learning Winter '17

Faster R-CNN (based on ResNets)

K. He, X. Zhang, S. Ren, J. Sun, [Deep Residual Learning for Image Recognition](#), CVPR 2016.

B. Leibe

43

Machine Learning Winter '17

YOLO

J. Redmon, S. Divvala, R. Girshick, A. Farhadi, [You Only Look Once: Unified, Real-Time Object Detection](#), CVPR 2016.

44

Machine Learning Winter '17

Object Detection Performance

PASCAL VOC

Slide credit: Ross Girshick

B. Leibe

45

Machine Learning Winter '17

Semantic Image Segmentation

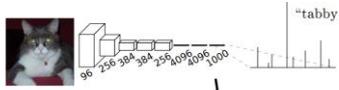
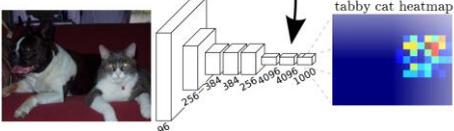
- Perform pixel-wise prediction task
 - Usually done using [Fully Convolutional Networks \(FCNs\)](#)
 - All operations formulated as convolutions
 - Advantage: can process arbitrarily sized images

Image source: Long, Shelhamer, Darrell

46

RWTH AACHEN UNIVERSITY

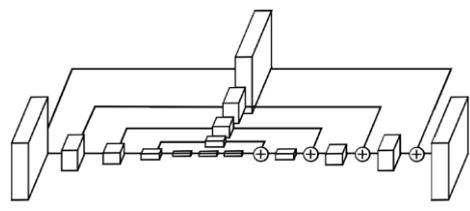
CNNs vs. FCNs

- CNN
 
- FCN
 
- Intuition
 - Think of FCNs as performing a sliding-window classification, producing a heatmap of output scores for each class

47
Image source: Loro, Shelhamer, Darabi

RWTH AACHEN UNIVERSITY

Semantic Image Segmentation



- Encoder-Decoder Architecture
 - Problem: FCN output has low resolution
 - Solution: perform upsampling to get back to desired resolution
 - Use skip connections to preserve higher-resolution information

48
Image source: Newell et al.

RWTH AACHEN UNIVERSITY

Semantic Segmentation

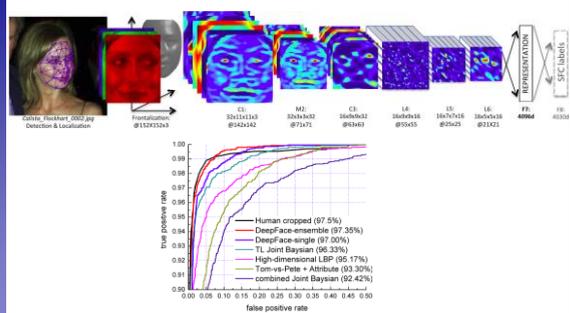


- Current state-of-the-art
 - Based on an extension of ResNets

(Pohlen, Hermans, Mathias, Leibe, CVPR 2017)

RWTH AACHEN UNIVERSITY

Other Tasks: Face Identification



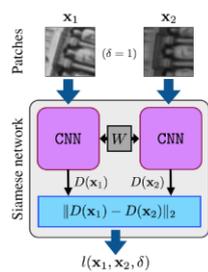
Y. Taigman, M. Yang, M. Ranzato, L. Wolf, **DeepFace: Closing the Gap to Human-Level Performance in Face Verification**, CVPR 2014

50
Slide credit: Svetlana Lazebnik

RWTH AACHEN UNIVERSITY

Learning Similarity Functions

- Siamese Network
 - Present the two stimuli to two identical copies of a network (with shared parameters)
 - Train them to output similar values if the inputs are (semantically) similar.
- Used for many matching tasks
 - Face identification
 - Stereo estimation
 - Optical flow
 - ...



B. Leibe

RWTH AACHEN UNIVERSITY

Extension: Triplet Loss Networks

- Learning a discriminative embedding
 - Present the network with triplets of examples
 - Negative
 - Anchor
 - Positive
 - Apply triplet loss to learn an embedding $f(\cdot)$ that groups the positive example closer to the anchor than the negative one.

$$\|f(x_i^a) - f(x_i^p)\|_2^2 < \|f(x_i^a) - f(x_i^n)\|_2^2$$


⇒ Used with great success in Google's FaceNet face identification

B. Leibe

References and Further Reading

- ResNets
 - K. He, X. Zhang, S. Ren, J. Sun, [Deep Residual Learning for Image Recognition](#), CVPR 2016.
 - A. Veit, M. Wilber, S. Belongie, [Residual Networks Behave Like Ensembles of Relatively Shallow Networks](#), NIPS 2016.

References: Computer Vision Tasks

- Object Detection
 - R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation, CVPR 2014.
 - S. Ren, K. He, R. Girshick, J. Sun, Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, NIPS 2015.
 - J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You Only Look Once: Unified Real-Time Object Detection, CVPR 2016.
 - W. Liu, D. Anguelov, [D. Erhan](#), [C. Szegedy](#), S. Reed, C-Y. Fu, A.C. Berg, SSD: Single Shot Multi Box Detector, ECCV 2016.

References: Computer Vision Tasks

- Semantic Segmentation
 - J. Long, E. Shelhamer, T. Darrell, Fully Convolutional Networks for Semantic Segmentation, CVPR 2015.
 - H. Zhao, J. Shi, X. Qi, X. Wang, J. Jia, Pyramid Scene Parsing Network, arXiv 1612.01105, 2016.