

RWTH AACHEN
UNIVERSITY

Machine Learning – Lecture 13

Convolutional Neural Networks

10.12.2018

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

Machine Learning Winter '18

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: Tricks of the Trade
- Convolutional Neural Networks
 - Neural Networks for Computer Vision
 - Convolutional Layers
 - Pooling Layers
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet

B. Leibe

RWTH AACHEN
UNIVERSITY

Recap: Reducing the Learning Rate

- Final improvement step after convergence is reached
 - Reduce learning rate by a factor of 10.
 - Continue training for a few epochs.
 - Do this 1-3 times, then stop training.
- Effect
 - Turning down the learning rate will reduce the random fluctuations in the error due to different gradients on different minibatches.
- *Be careful: Do not turn down the learning rate too soon!*
 - Further progress will be much slower/impossible after that.

B. Leibe

RWTH AACHEN
UNIVERSITY

Recap: Data Augmentation

- Effect
 - Much larger training set
 - Robustness against expected variations
- During testing
 - When cropping was used during training, need to again apply crops to get same image size.
 - Beneficial to also apply flipping during test.
 - Applying several ColorPCA variations can bring another ~1% improvement, but at a significantly increased runtime.

Augmented training data
(from one original image)

B. Leibe

RWTH AACHEN
UNIVERSITY

Recap: Normalizing the Inputs

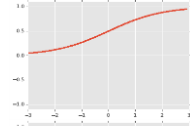
- Convergence is fastest if
 - The mean of each input variable over the training set is zero.
 - The inputs are scaled such that all have the same covariance.
 - Input variables are uncorrelated if possible.
- Advisable normalization steps (for MLPs only, not for CNNs)
 - Normalize all inputs that an input unit sees to zero-mean, unit covariance.
 - If possible, try to decorrelate them using PCA (also known as Karhunen-Loeve expansion).

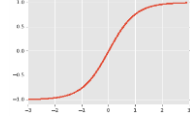
B. Leibe

Machine Learning Winter '18

Recap: Commonly Used Nonlinearities

- Sigmoid

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

- Hyperbolic tangent

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

- Softmax

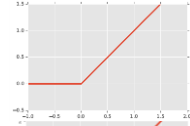
$$g(\mathbf{a}) = \frac{\exp\{-a_i\}}{\sum_j \exp\{-a_j\}}$$


B. Leibe 7

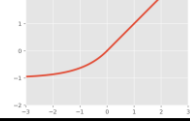
Machine Learning Winter '18

Recap: Commonly Used Nonlinearities (2)

- Rectified linear unit (ReLU)

$$g(a) = \max\{0, a\}$$

- Leaky ReLU

$$g(a) = \max\{\beta a, a\} \quad \beta \in [0.01, 0.3]$$
 - Avoids stuck-at-zero units
 - Weaker offset bias
- ELU

$$g(a) = \begin{cases} a, & a \geq 0 \\ e^a - 1, & a < 0 \end{cases}$$
 - No offset bias anymore
 - BUT: need to store activations

B. Leibe 8

Machine Learning Winter '18

Recap: Glorot Initialization [Glorot & Bengio, '10]

- Variance of neuron activations
 - Suppose we have an input X with n components and a linear neuron with random weights W that spits out a number Y .
 - We want the variance of the input and output of a unit to be the same, therefore $n \text{Var}(W_i)$ should be 1. This means

$$\text{Var}(W_i) = \frac{1}{n} = \frac{1}{n_{\text{in}}}$$
 - Or for the backpropagated gradient

$$\text{Var}(W_i) = \frac{1}{n_{\text{out}}}$$
 - As a compromise, Glorot & Bengio propose to use

$$\text{Var}(W) = \frac{2}{n_{\text{in}} + n_{\text{out}}}$$

⇒ Randomly sample the weights with this variance. That's it.

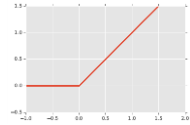
B. Leibe 9

Machine Learning Winter '18

Recap: He Initialization [He et al., '15]

- Extension of Glorot Initialization to ReLU units
 - Use Rectified Linear Units (ReLU)

$$g(a) = \max\{0, a\}$$
 - Effect: gradient is propagated with a constant factor

$$\frac{\partial g(a)}{\partial a} = \begin{cases} 1, & a > 0 \\ 0, & \text{else} \end{cases}$$

- Same basic idea: Output should have the input variance
 - However, the Glorot derivation was based on tanh units, linearity assumption around zero does not hold for ReLU.
 - He et al. made the derivations, proposed to use instead

$$\text{Var}(W) = \frac{2}{n_{\text{in}}}$$

B. Leibe 10

Machine Learning Winter '18

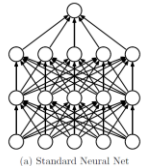
Recap: Batch Normalization [Ioffe & Szegedy '14]

- Motivation
 - Optimization works best if all inputs of a layer are normalized.
- Idea
 - Introduce intermediate layer that centers the activations of the previous layer per minibatch.
 - I.e., perform transformations on all activations and undo those transformations when backpropagating gradients
 - **Complication:** centering + normalization also needs to be done at test time, but minibatches are no longer available at that point.
 - Learn the normalization parameters to compensate for the expected bias of the previous layer (usually a simple moving average)
- Effect
 - Much improved convergence (but parameter values are important!)
 - Widely used in practice

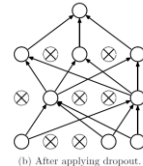
B. Leibe 11

Machine Learning Winter '18

Recap: Dropout [Srivastava, Hinton '12]



(a) Standard Neural Net



(b) After applying dropout.

- Idea
 - Randomly switch off units during training.
 - Change network architecture for each data point, effectively training many different variants of the network.
 - When applying the trained network, multiply activations with the probability that the unit was set to zero.

⇒ Greatly improved performance

B. Leibe 12

Machine Learning Winter '18

Topics of This Lecture

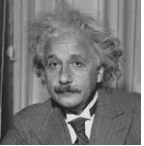
- Recap: Tricks of the Trade
- Convolutional Neural Networks
 - Neural Networks for Computer Vision
 - Convolutional Layers
 - Pooling Layers
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet

B. Leibe 13

Machine Learning Winter '18

Neural Networks for Computer Vision

- How should we approach vision problems?


→ Face Y/N?
- Architectural considerations
 - Input is 2D ⇒ 2D layers of units
 - No pre-segmentation ⇒ Need robustness to misalignments
 - Vision is hierarchical ⇒ Hierarchical multi-layered structure
 - Vision is difficult ⇒ Network should be deep

B. Leibe 14

Machine Learning Winter '18

Why Hierarchical Multi-Layered Models?

- Motivation 1: Visual scenes are hierarchically organized

Object

↑

Object parts

↑

Primitive features

↑

Input image

Face

↑

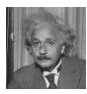
Eyes, nose, ...

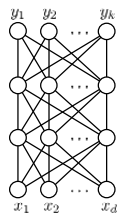
↑

Oriented edges

↑

Face image





Slide adapted from Richard Turner. B. Leibe 15

Machine Learning Winter '18

Why Hierarchical Multi-Layered Models?

- Motivation 2: *Biological vision* is hierarchical, too

Object

↑

Object parts

↑

Primitive features

↑

Input image

Face

↑

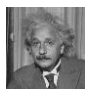
Eyes, nose, ...

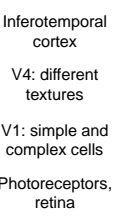
↑

Oriented edges

↑

Face image






Inferotemporal cortex

V4: different textures

V1: simple and complex cells

Photoreceptors, retina



Slide adapted from Richard Turner. B. Leibe 16

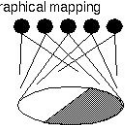
Machine Learning Winter '18

Hubel/Wiesel Architecture

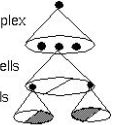
- D. Hubel, T. Wiesel (1959, 1962, Nobel Prize 1981)
 - Visual cortex consists of a hierarchy of *simple*, *complex*, and *hyper-complex* cells

Hubel & Wiesel

topographical mapping



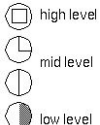
featural hierarchy



high level

mid level

low level



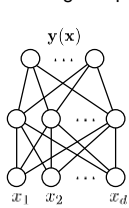
Slide credit: Svetlana Lazebnik, Rob Fergus. B. Leibe 18

Machine Learning Winter '18

Why Hierarchical Multi-Layered Models?

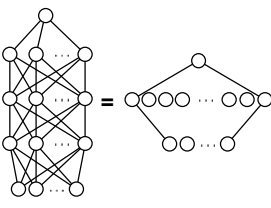
- Motivation 3: Shallow architectures are inefficient at representing complex functions

$y(x)$



$x_1 \ x_2 \ \dots \ x_d$

An MLP with 1 hidden layer can implement *any* function (universal approximator)



However, if the function is deep, a very large hidden layer may be required.

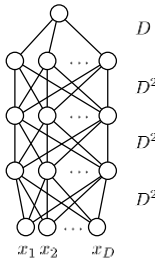
Slide adapted from Richard Turner. B. Leibe 19

What's Wrong With Standard Neural Networks?

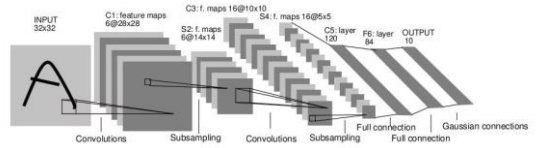
- Complexity analysis
 - How many parameters does this network have?

$$|\theta| = 3D^2 + D$$
 - For a small 32×32 image

$$|\theta| = 3 \cdot 32^4 + 32^2 \approx 3 \cdot 10^6$$
- Consequences
 - Hard to train
 - Need to initialize carefully
 - *Convolutional nets reduce the number of parameters!*



Convolutional Neural Networks (CNN, ConvNet)

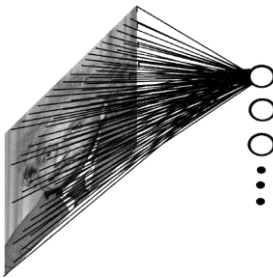


- Neural network with specialized connectivity structure
 - Stack multiple stages of feature extractors
 - Higher stages compute more global, more invariant features
 - Classification layer at the end

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.

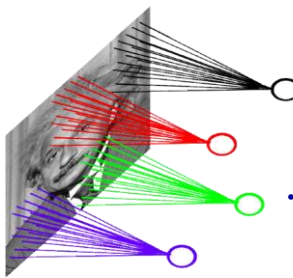
Convolutional Networks: Intuition

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



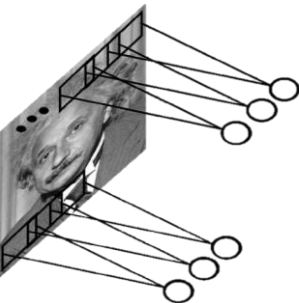
Convolutional Networks: Intuition

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



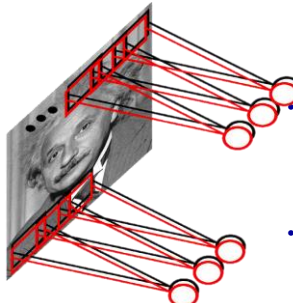
Convolutional Networks: Intuition

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



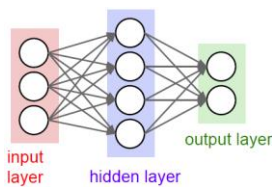
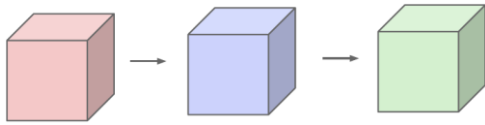
Convolutional Networks: Intuition

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



RWTH AACHEN UNIVERSITY

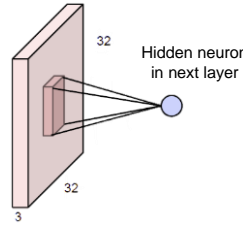
Important Conceptual Shift

- Before
 
- Now:
 

Machine Learning Winter '18 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 26

RWTH AACHEN UNIVERSITY

Convolution Layers



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

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

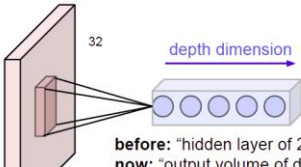
Now: Local connectivity
One neuron connects to, e.g., $5 \times 5 \times 3$ region.
 \Rightarrow Only $5 \times 5 \times 3$ shared weights.

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

Machine Learning Winter '18 | Slide adapted from FeiFei Li, Andrei Karpathy | B. Leibe | 27

RWTH AACHEN UNIVERSITY

Convolution Layers



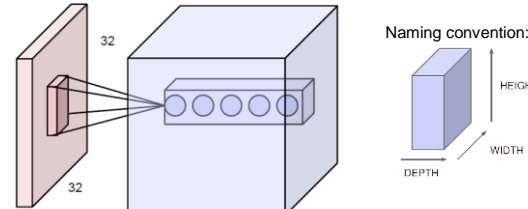
before: "hidden layer of 200 neurons"
now: "output volume of depth 200"

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

Machine Learning Winter '18 | Slide adapted from FeiFei Li, Andrei Karpathy | B. Leibe | 28

RWTH AACHEN UNIVERSITY

Convolution Layers



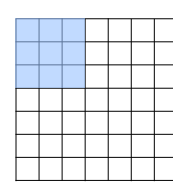
Naming convention:

- 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.

Machine Learning Winter '18 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 29

RWTH AACHEN UNIVERSITY

Convolution Layers



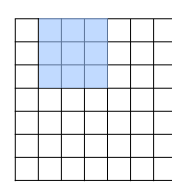
Example:
 7×7 input
assume 3×3 connectivity
stride 1

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

Machine Learning Winter '18 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 31

RWTH AACHEN UNIVERSITY

Convolution Layers



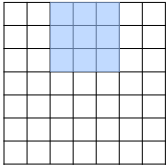
Example:
 7×7 input
assume 3×3 connectivity
stride 1

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

Machine Learning Winter '18 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 32

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**.

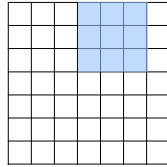
Machine Learning Winter '18

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

33

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**.

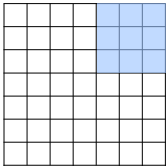
Machine Learning Winter '18

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

34

RWTH AACHEN UNIVERSITY

Convolution Layers



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

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

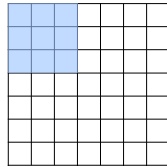
Machine Learning Winter '18

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

35

RWTH AACHEN UNIVERSITY

Convolution Layers



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

What about stride 2?

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

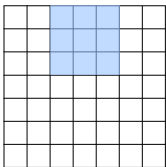
Machine Learning Winter '18

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

36

RWTH AACHEN UNIVERSITY

Convolution Layers



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

What about stride 2?

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

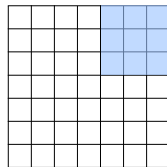
Machine Learning Winter '18

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

37

RWTH AACHEN UNIVERSITY

Convolution Layers



Example:
 7×7 input
 assume 3×3 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**.

Machine Learning Winter '18

Slide credit: FeiFei Li, Andrei Karpathy B. Leibe

38

RWTH AACHEN UNIVERSITY

Convolution Layers

0	0	0	0	0			
0							
0							
0							
0							

Example:
 7×7 input
 assume 3×3 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.

Machine Learning Winter '18 | Slide credit: FeiFei Li, Andrei Karpathy | B. Leibe | 39

RWTH AACHEN UNIVERSITY

Activation Maps of Convolutional Filters

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

5x5 filters

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

Activation maps

Machine Learning Winter '18 | Slide adapted from FeiFei Li, Andrei Karpathy | B. Leibe | 40

RWTH AACHEN UNIVERSITY

Effect of Multiple Convolution Layers

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

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

Machine Learning Winter '18 | Slide credit: Yann LeCun | B. Leibe | 41

RWTH AACHEN UNIVERSITY

Convolutional Networks: Intuition

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

Machine Learning Winter '18 | Slide adapted from Marc'Aurelio Ranzato | B. Leibe | Image source: Yann LeCun | 42

RWTH AACHEN UNIVERSITY

Convolutional Networks: Intuition

- 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.

Machine Learning Winter '18 | Slide adapted from Marc'Aurelio Ranzato | B. Leibe | Image source: Yann LeCun | 43

RWTH AACHEN UNIVERSITY

Max Pooling

Single depth slice

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

max pool with 2×2 filters and stride 2

6	8
3	4

Effect:

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

Machine Learning Winter '18 | Slide adapted from FeiFei Li, Andrei Karpathy | B. Leibe | 44

RWTH AACHEN UNIVERSITY

Max Pooling

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

x
y

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

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

45

RWTH AACHEN UNIVERSITY

CNNs: Implication for Back-Propagation

- Convolutional layers
 - Filter weights are shared between locations
 - Gradients are added for each filter location.

Machine Learning Winter '18 B. Leibe

46

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Recap: Tricks of the Trade
- Convolutional Neural Networks
 - Neural Networks for Computer Vision
 - Convolutional Layers
 - Pooling Layers
- CNN Architectures
 - LeNet
 - AlexNet
 - VGGNet
 - GoogLeNet

Machine Learning Winter '18 B. Leibe

47

RWTH AACHEN UNIVERSITY

CNN Architectures: LeNet (1998)

- Early convolutional architecture
 - 2 Convolutional layers, 2 pooling layers
 - Fully-connected NN layers for classification
 - Successfully used for handwritten digit recognition (MNIST)

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

48

RWTH AACHEN UNIVERSITY

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]

Machine Learning Winter '18 B. Leibe

49

RWTH AACHEN UNIVERSITY

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.E. Hinton, NIPS 2012

Machine Learning Winter '18 B. Leibe

50

Machine Learning Winter '18

More Finegrained Classes

PASCAL

birds

bird

ILSVRC

flamingo cock ruffed grouse quail partridge ...

cats

cat Egyptian cat Persian cat Siamese cat tabby lynx ...

dogs

dog dalmatian keeshond miniature schnauzer standard schnauzer giant schnauzer ...

B. Leibe Image source: O. Russakovsky et al.

64

Machine Learning Winter '18

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 Image source: A. Krizhevsky

65

Machine Learning Winter '18

References and Further Reading

- LeNet
 - 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.
- 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

66

Machine Learning Winter '18

References and Further Reading

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

B. Leibe

67