

Word Spotting: From Bag-of-Features to Deep Learning

— **SSDA 2017 Tutorial, Jaipur, India** —

Gernot A. Fink & Sebastian Sudholt

January, 2017

- ▶ Introduction
- ▶ Bag-of-Features: Fundamentals
- ▶ Learning Document Image Representations
- ▶ Learning Word Spotting Models
- ▶ Discussion
- ▶ Deep Learning Fundamentals
- ▶ Deep Learning for Word Spotting
- ▶ Summary

Introduction: Automatic Reading Systems

State of Automatic Reading:

- ▶ One of the earliest application fields studied in computer science
- ▶ So-called OCR achieves high-quality results for machine-printed text in well-defined settings.
- ▶ Online handwriting recognition again gaining popularity
- ▶ Offline handwriting recognition: Remarkable results, but still an open research problem

General Methodology:

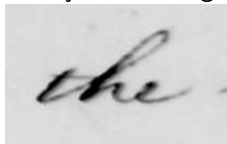
Statistical sequence models (e.g. Hidden-Markov Models) that are trained from *extensive* amounts of example data

Introduction: Why Word Spotting?

What if automatic transcription of handwriting is no longer feasible?

Alternative: Retrieval of individual words rather than transcription (“query-by-example”)

Query word image



Document image

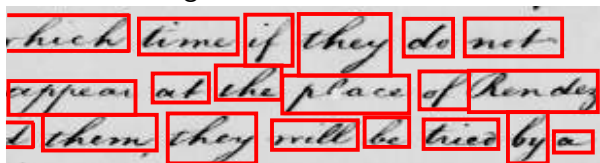
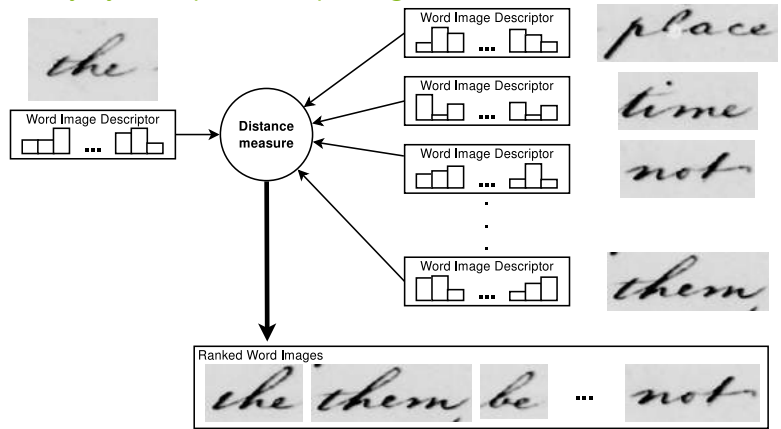


Image of Feldpost postcard from *Private Collection Dr. Britta Bley*, Dortmund, Germany Images from *The George Washington Papers at the Library of Congress, 1741-1799*

Introduction: Basic Methodology

Query-by-example word spotting



Based on [Rath & Manmatha, IJRAR'07]

Word Spotting: Fundamentals

Core Methodology:

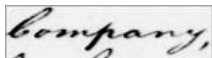
- ▶ Specialized image retrieval
- ▶ Important ingredient: Image matching procedure
- ▶ Frequently required: Pre-segmentation (words / lines)

Taxonomy:

- ▶ *Segmentation-based*
- ▶ *Segmentation-free*, i.e., segmentation problem covered during retrieval
- ▶ *Query-by-Example*, i.e., word image directly used as query
- ▶ *Query-by-String*, i.e., query model derived from textual query (“string”)

Word Spotting Tasks

Query by Example



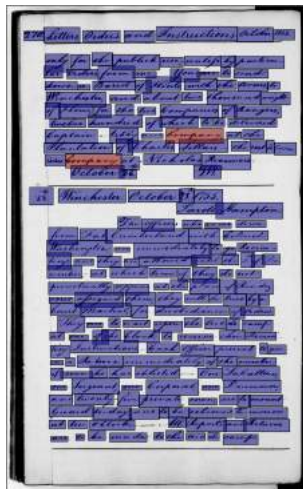
company,



Query by String



Company



Word-Spotting: Milestones

Manmatha *et al.* 1996: First influential work

(Binarization, Alignment, XOR distance)

Rath & Manmatha 2003: DTW matching

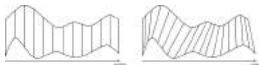
(Normalization, profile features)

Rusiñol *et al.* 2011: First influential work
 using BoF, first with Spatial Pyramid

(SIFT, BoF, Spatial Pyramid, LSI,
 segmentation-free decoding)

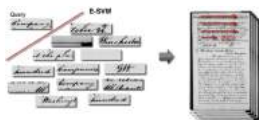
Almazan *et al.* 2014: HOG features

(“Exemplar SVM”, query expansion)



(a) naive alignment after resampling. (b) alignment with DTW.

Company	Company	Company	Company
English	but English will	an English man	English Ho
است	ی است	ی است	ی است
	ی است	ی است	ی است



Overview

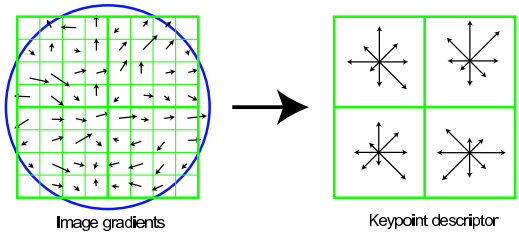
- ▶ Introduction
- ▶ **Bag-of-Features: Fundamentals**
- ▶ Learning Document Image Representations
- ▶ Learning Word Spotting Models
- ▶ Discussion
- ▶ Deep Learning Fundamentals
- ▶ Deep Learning for Word Spotting
- ▶ Summary

Local Image Descriptors

Fundamentals:

- ▶ Local gradient statistics (i.e. implicit description of object contours)
- ▶ Grouping of multiple such local statistics in a certain neighborhood and *normalization* (coarse description of structural properties)

Example: SIFT descriptor [Lowe, 2004]



(Source: [Lowe, IJCV 2004])

Note: Ex. several similar descriptors, e.g., HOG, SURF

Statistical Image Modeling

Representation: Bag-of-Features Models (BoF)

- ▶ Originally proposed as Bag-of-Words models for representing texts

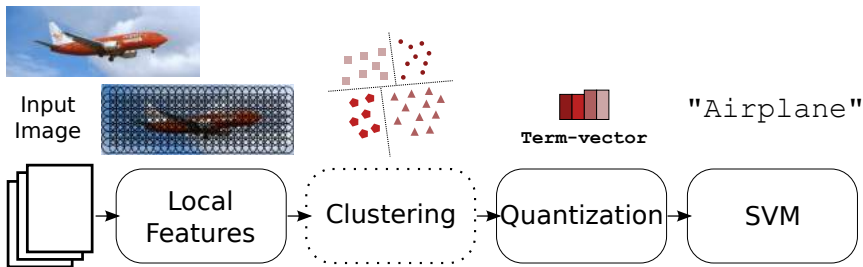
BoF-Approach: (cf. e.g. [O'Hara & Draper 2011])

0. Extract local image features (usually in a grid)
1. Compute *visual vocabulary* by quantizing local image features
2. For given image, compute histogram of quantized features (i.e. orderless “bag” of features)
3. Any classification / matching technique can be applied to BoF representations

S. O'Hara & B. A. Draper: [Introduction to the Bag of Features Paradigm for Image Classification and Retrieval](#), Computing Research Repository, arXiv:1101.3354v1, 2011.

Statistical Image Modeling II

BoF-Approach: Processing Pipeline

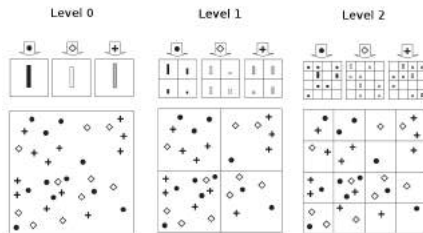


Main drawback: Spatial information is lost!

Statistical Image Modeling III

Goal: Compensate loss of spatial information

Method: Spatial-pyramid models



Lazebnik, S., Schmid, C., Ponce, J.: *Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories*, IEEE Conference on Computer Vision and Pattern Recognition, 2006.

Disadvantage: Considerably increased model complexity

Overview

- ▶ Introduction
- ▶ Bag-of-Features: Fundamentals
- ▶ **Learning Document Image Representations**
- ▶ Learning Word Spotting Models
- ▶ Discussion
- ▶ Deep Learning Fundamentals
- ▶ Deep Learning for Word Spotting
- ▶ Summary

Bag-of-Features Models for Word Spotting

Basic Methodology:

Image Features:

Gradient-based descriptors, e.g., SIFT, HOG

Feature Extraction:

Dense grid, i.e., no keypoint detection involved

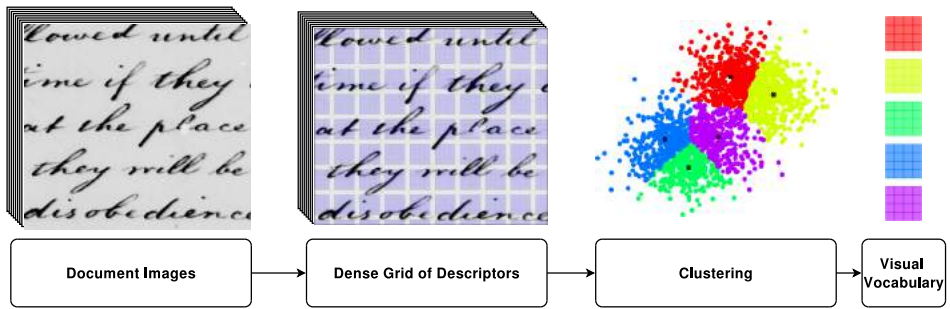
Visual Vocabulary:

- ▶ Quantization of descriptors (as “usual”)
- ▶ **Special:** Large vocabularies (i.e. 2K to 4K) in order to capture appearance of “individuals” precisely

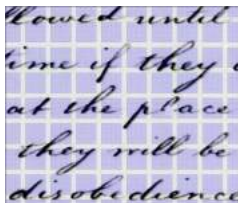
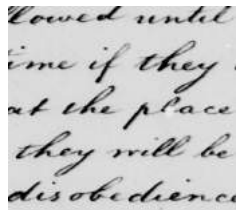
Query Model:

- ▶ Segmentation- or patch-based processing
- ▶ 1D spatial pyramid for improved spatial modeling

Bag-of-Features Models: Visual Vocabulary



Bag-of-Features Models: Term Vector



Document Images

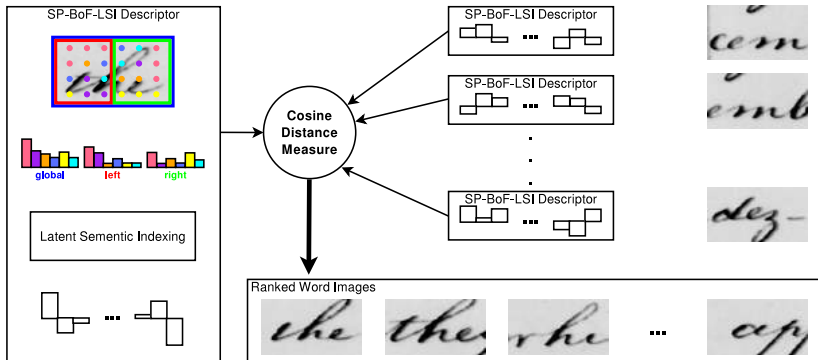
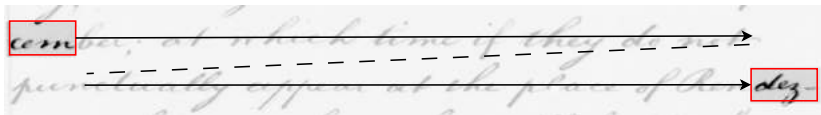
Dense Grid of Descriptors

Quantization

Visual Vocabulary

Bag-of-Features
Term Vector

Segmentation-free Word Spotting with Bag-of-Features

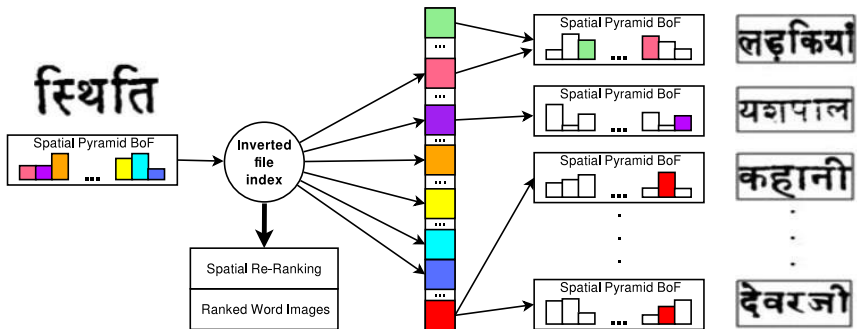


Rusiñol, M., and Aldavert, D., and Toledo, R., and Lladós, J.: [Browsing heterogeneous document collections by a segmentation-free word spotting method](#), Int. Conf. on Document Analysis and Recognition, Beijing, pp. 63-67, 2011.

Fink, Sudholt

Word Spotting: BoF to Deep Learning

Segmentation-based Word Spotting with Bag-of-Features



Shekhar, R., and Jawahar, C.: [Word image retrieval using bag of visual words](#), Int. Workshop on Document Analysis Systems, pp. 297-301, 2012.

Overview

- ▶ Introduction
- ▶ Bag-of-Features: Fundamentals
- ▶ Learning Document Image Representations
- ▶ **Learning Word Spotting Models**
- ▶ Discussion
- ▶ Deep Learning Fundamentals
- ▶ Deep Learning for Word Spotting
- ▶ Summary

Learning BoF-Based Word Spotting Models

What we have so far:

- ▶ Expert-designed image descriptors
- ▶ Automatically learned document image features
(visual vocabulary → histograms of quantized descriptors)
- ▶ Matching: Nearest neighbor

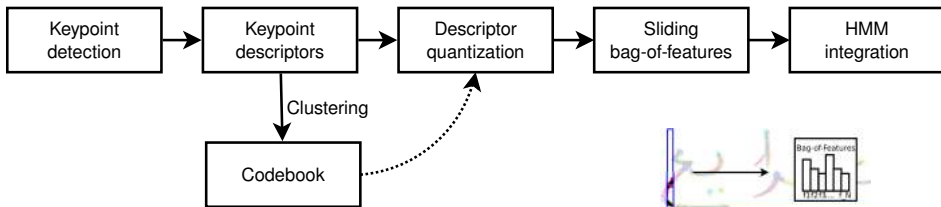
What we want in addition:

- ▶ More powerful image matching
- ▶ Categorization capabilities
(i.e. links between image appearance and textual representation)

What we need: Learned classification models

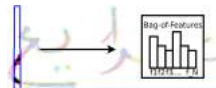
Bag-of-Features HMMs

- ▶ *Extension* of HMMs towards learned feature representation
- ▶ *Extension* of BoF models towards fine-grained script representation



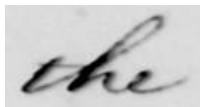
$$b_j(\mathbf{f}) = \sum_{k=1}^{|\mathcal{V}|} c_{jk} f_k$$

with \mathbf{f} : term vector
 \mathcal{V} : vis. voc.

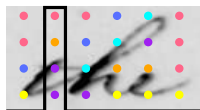


Rothacker, L., Vajda, S., Fink, G. A.: *Bag-of-Features Representations for Offline Handwriting Recognition Applied to Arabic Script*, In Proc. ICFHR, Bari, 2012.

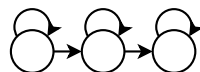
Bag-of-Features HMMs for QbE Word Spotting



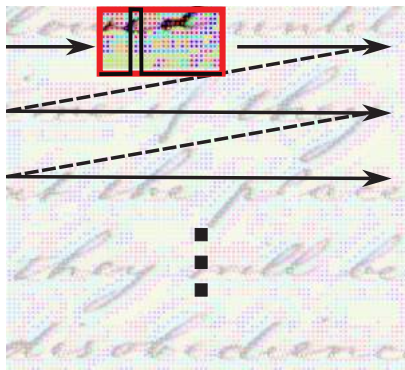
Query Word



Bag-of-Features Sequence (sliding window)



Bag-of-Features Hidden Markov Model



Approach can be sped-up by intelligent patch pre-filtering!

Rothacker, L., Rusinol, M., Fink, G. A.: *Bag-of-Features HMMs for Segmentation-Free Word Spotting in Handwritten Documents*, In Proc. Int. Conf. on Document Analysis and Recognition, Washington DC, USA, 2013.

Rothacker, L., Rusinol, M., Lladós, J., Fink, G. A.: *A Two-Stage Approach to Segmentation-Free Query-by-Example Word Spotting*, manuscript cultures, 1(7), pages 47-57, 2014.

Bag-of-Features HMMs for QbS Word Spotting

Training set

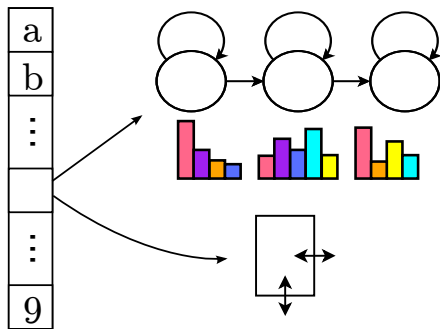


c a p t a i n

⋮



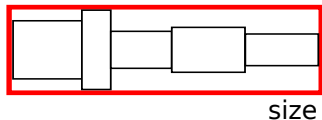
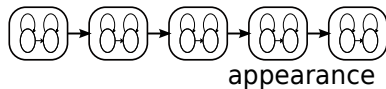
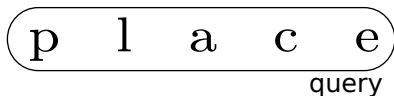
t w e l v e



- ▶ Appearance model: Bag-of-Features HMMs (per character)
- ▶ Spatial size model: Width & height estimates

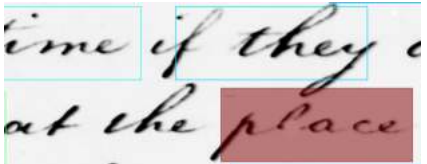
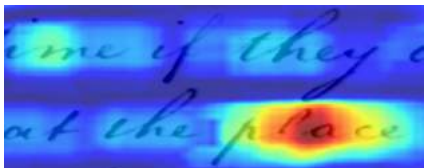
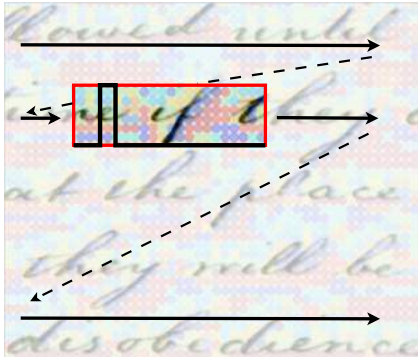
Rothacker, L., Fink, G. A.: *Segmentation-Free Query-by-String Word Spotting with Bag-of-Features HMMs*, Int. Conf. on Document Analysis and Recognition, Nancy, France, 2015.

Bag-of-Features HMMs for QbS Word Spotting II



Rothacker, L., Fink, G. A.: *Segmentation-Free Query-by-String Word Spotting with Bag-of-Features HMMs*, Int. Conf. on Document Analysis and Recognition, Nancy, France, 2015.

Bag-of-Features HMMs for QbS Word Spotting II



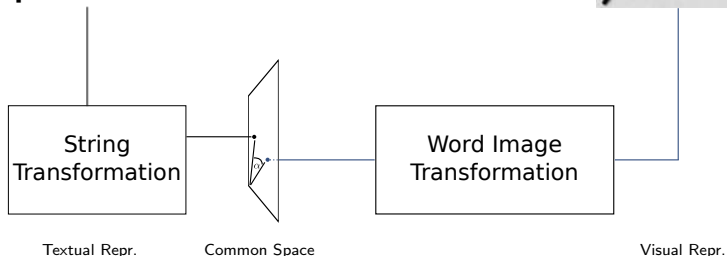
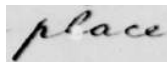
Rothacker, L., Fink, G. A.: *Segmentation-Free Query-by-String Word Spotting with Bag-of-Features HMMs*, Int. Conf. on Document Analysis and Recognition, Nancy, France, 2015.

Subspace Representations for Word Spotting

Idea: Project both textual and visual representation into a *common space*

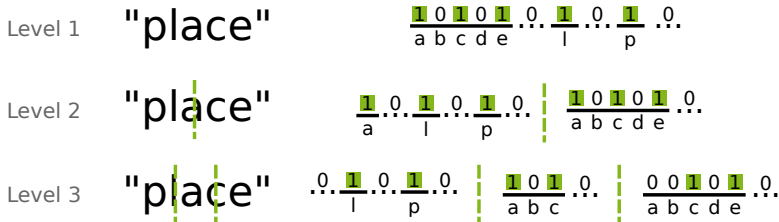
Benefits: QbE and QbS are now a simple nearest neighbor search

"place"



J. Almazán, A. Gordo, A. Fornés and E. Valveny: [Word Spotting and Recognition with Embedded Attributes](#), IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 36, no. 12, pp. 2552-2566, 2014.

Pyramidal Histogram of Characters (PHOC)

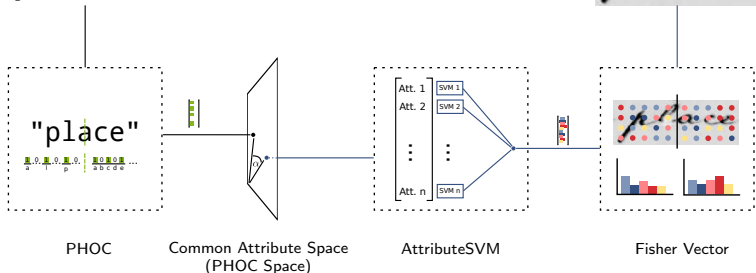


- ▶ Concatenate histograms for all levels to form PHOC
- ▶ Levels used by Almazán *et al.*: 2,3,4 and 5
- ▶ 26 Characters + 10 Digits
- ▶ $\text{PHOC} \in \{0, 1\}^{604}$

J. Almazán, A. Gordo, A. Fornés and E. Valveny: [Word Spotting and Recognition with Embedded Attributes](#), IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 36, no. 12, pp. 2552-2566, 2014.

Learning the PHOC representation

"place"



- ▶ AttributeSVM: ensemble of SVMs
- ▶ each SVM predicts the presence or absence of one element of the PHOC

J. Almazán, A. Gordo, A. Fornés and E. Valveny: [Word Spotting and Recognition with Embedded Attributes](#), IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 36, no. 12, pp. 2552-2566, 2014.

Overview

- ▶ Introduction
- ▶ Bag-of-Features: Fundamentals
- ▶ Learning Document Image Representations
- ▶ Learning Word Spotting Models
- ▶ **Discussion**
- ▶ Deep Learning Fundamentals
- ▶ Deep Learning for Word Spotting
- ▶ Summary

Pros and Cons of Methods so Far

Basic BoF-Based Models:

- ✓ (partly) Learned image features / representations
- ✓ No annotations required (unsupervised)
- ⚡ Purely image-based (no categorial information)

Advanced BoF-Based Models (BoF-HMMs, AttributeSVMs):

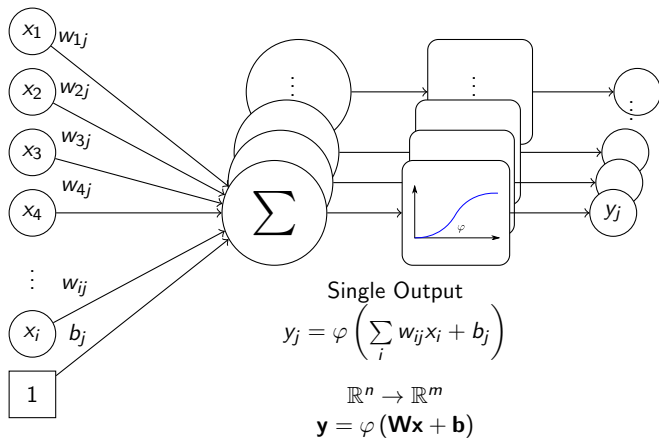
- ✓ Learned (categorial) models (supervised)
- ✓ Common representation of image appearance and categorial information (embedded attributes)
- ⚡ Features and model learned separately

Desired: Integrated framework with overall / end-to-end optimization

Overview

- ▶ Introduction
- ▶ Bag-of-Features: Fundamentals
- ▶ Learning Document Image Representations
- ▶ Learning Word Spotting Models
- ▶ Discussion
- ▶ **Deep Learning Fundamentals**
- ▶ Deep Learning for Word Spotting
- ▶ Summary

The Perceptron

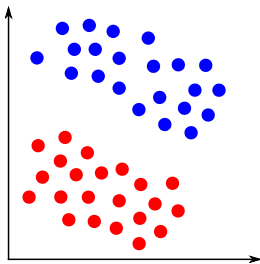


F. Rosenblatt: [The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain](#), Psychological Review, 65(6), 1958.

Capabilities of the Perceptron

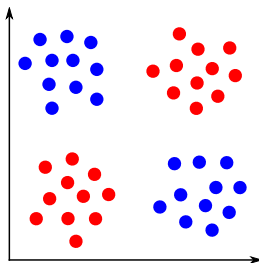
What a Perceptron can do:

Classify two linearly separable classes



What a Perceptron can't do:

Classify two non-linearly separable classes (XOR-Problem)

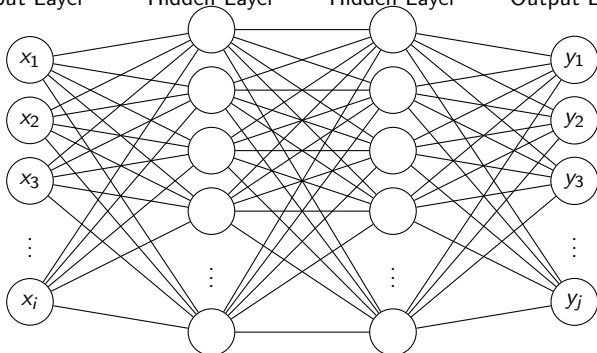


Solution: Stack layers of Perceptrons

⇒ Multi Layer Perceptron

Multi Layer Perceptron (MLP)

Input Layer Hidden Layer Hidden Layer Output Layer



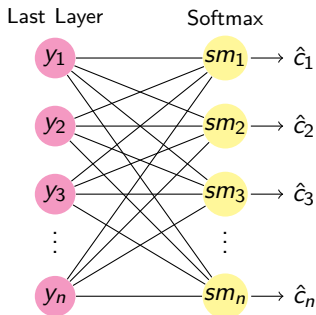
$$\mathbf{y} = \mathbf{f}^L \left(\mathbf{f}^{L-1} (\dots \mathbf{f}^2 (\mathbf{f}^1 (\mathbf{x})) \dots) \right)$$

here: $\mathbf{y} = \mathbf{W}_{\text{out}} \cdot \varphi (\mathbf{W}_{h2} \cdot \varphi (\mathbf{W}_{h1} \mathbf{x} + b_{h1}) + b_{h2}) + b_{\text{out}}$

Note: Activation function is typically left out for last layer

Classifying with MLPs

- ▶ For classification, the output of the MLP is forwarded through a Softmax Function: $sm_i(\mathbf{y}) = \frac{e^{y_i}}{\sum_j e^{y_j}}$
- ▶ Softmax can be seen as an additional layer in the MLP
- ▶ sm_i is pseudo-probability for class c_i
- ▶ Predicted class: $\hat{c} = \max_i sm_i$

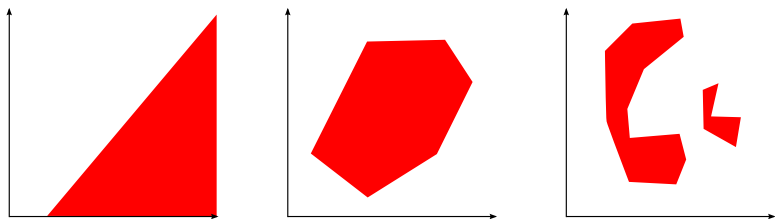


What an MLP Can Do!

... approximate any function (even with only 2 layers!)

[Hornik *et al.* 1989]

Interpretation with 3 layers (2 hidden, 1 output):



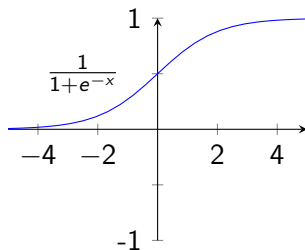
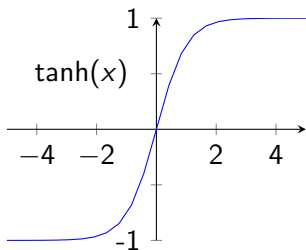
1. Layer: Halfspaces
2. Layer: Convex polyhedron
3. Layer: Multiple non-convex, non-connected polyhedra

A Word on Activation Functions

- ▶ Activation functions are crucial for MLP
- ▶ Without non-linearities, an MLP implements a linear transform:

$$\mathbf{y} = \mathbf{W}_L \mathbf{W}_{L-1} \dots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} = \mathbf{W}' \mathbf{x}$$

Classic Activation Functions: sigmoidal shape (“threshold-like”)



Training an MLP

How to determine weights such that desired function is performed?

Basic Idea: Compare (computed) output of MLP

$$\hat{\mathbf{y}} = \mathbf{f}^L \left(\mathbf{f}^{L-1} (\dots \mathbf{f}^2 (\mathbf{f}^1 (\mathbf{x}))) \right)$$

to *desired* (ideal) output \mathbf{y} and

update weights such that $\hat{\mathbf{y}}$ and \mathbf{y} become more similar.

Comparison requires *loss function* that evaluates similarity of $\hat{\mathbf{y}}$ and \mathbf{y} , e.g., Cross-Entropy (in comb. w. Softmax):

$$\epsilon = - \sum y_i \log \hat{y}_i$$

Update weights in the negative direction of the gradient of the loss:

$$w_{ij}^l \leftarrow w_{ij}^l - \beta \frac{\partial \epsilon}{\partial w_{ij}^l}$$

⇒ Training $\hat{=}$ Gradient Descent

Training an MLP: Error Backpropagation

How to compute gradient for network weights?

Final layer: straight forward, i.e., using $f_j^l = \phi(g_j^l)$:

$$\frac{\partial \epsilon}{\partial w_{ij}^L} = \frac{\partial \epsilon}{\partial f_j^L} \cdot \frac{\partial f_j^L}{\partial g_j^L} \cdot \frac{\partial g_j^L}{\partial w_{ij}^L}$$

Hidden layers: define gradient based on "local error" $\delta_j^k = \frac{\partial \epsilon}{\partial g_j^k}$:

$$\frac{\partial \epsilon}{\partial w_{ij}^k} = \frac{\partial \epsilon}{\partial \mathbf{f}^L} \cdot \frac{\partial \mathbf{f}^L}{\partial \mathbf{f}^{L-1}} \cdot \dots \cdot \frac{\partial \mathbf{f}^{k+1}}{\partial \mathbf{f}^k} \cdot \frac{\partial \mathbf{f}^k}{\partial w_{ij}^k} \frac{\partial \mathbf{f}^k}{\partial g_j^k} \cdot \frac{\partial g_j^k}{\partial w_{ij}^k} = \delta_j^k \cdot \frac{\partial g_j^k}{\partial w_{ij}^k}$$

Local error can be computed *backward* through the network:

$$\delta_j^l = \left\{ \sum w_{jk}^{l+1} d_k^{l+1} \right\} \cdot \frac{\partial f^{l+1}}{\partial g_j^{l+1}}$$

⇒ Error Backpropagation

Classifying Images with Neural Networks

Problem:

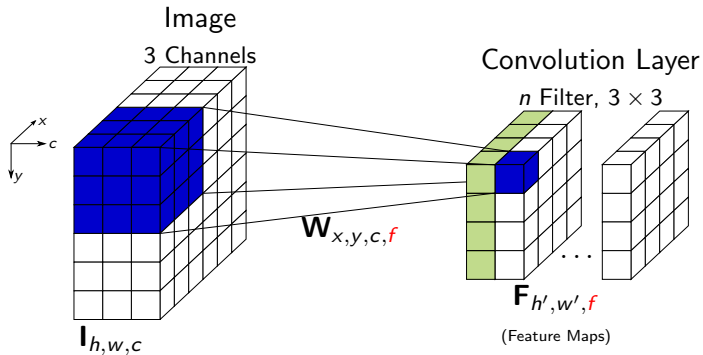
- ▶ Using MLPs for image classification is only possible for very small images (e.g. 28×28 pixels)
- ▶ Number of weights would explode for bigger images

Example: RGB Image of 224×224 pixels,
first hidden MLP layer has 768 neurons (small layer):
 $224 \cdot 224 \cdot 3 \cdot 768 \approx 10^8$ weights in the first layer (441 MB)

Solution: Don't use fully connected layer but rather apply a small number of weights at all possible locations in the image

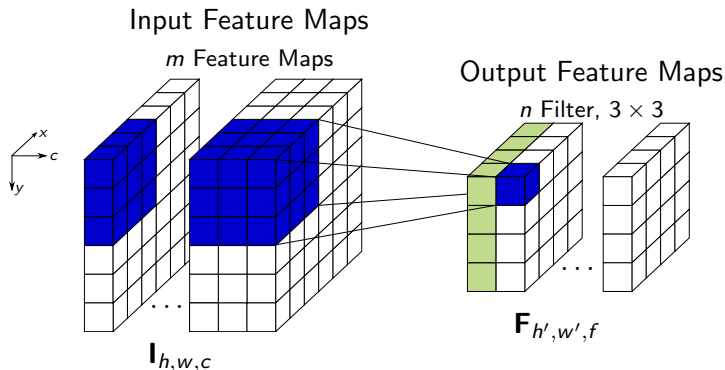
⇒ Convolutional Layer

Convolutional Layer

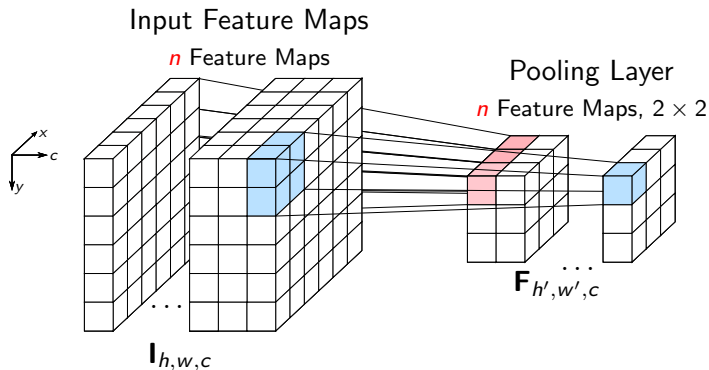


$$F_{x,y,f} = \varphi \left(\sum_{c=1}^K \sum_{i=1}^3 \sum_{j=1}^3 W_{i,j,c,f} \cdot I_{x+i,y+j,c} + b_f \right)$$

Cascade of Convolutional Layers

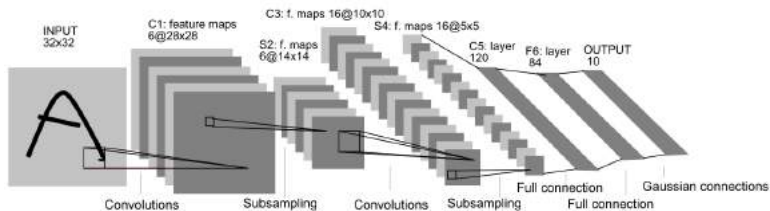


Pooling Layer



$$F_{x,y,f} = \max_{i,j} I_{x+i,y+j,f}$$

LeNet



(Source: [LeCun et al., 1990])

- ▶ LeNet predicts one of 10 character classes for a given input image
- ▶ Subsampling = Pooling Layer
- ▶ Gaussian Connections = FC Layer + Euclidean Loss

Y. LeCun, B. Boser, J.S. Denker, D. Henderson, R. E. Howard, W. Hubbard, L.D. Jackel: [Handwritten Digit Recognition with a Back-Propagation Network](#), Neural Information Processing Systems, pp. 396–404, 1990.

Deep Learning

In general: Deeper network architectures perform better than shallower ones for vision tasks

Important: Only empirical evidence (no theoretical proofs)

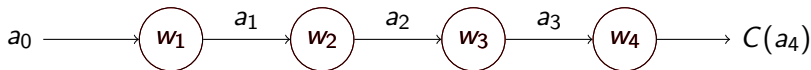
Technically: Deeper means more layers, not a deeper understanding

Even with high computation power and large datasets,
Deep Learning did not really pick up until 2012

Why? Vanishing Gradient Problem

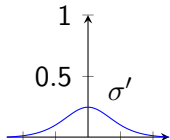
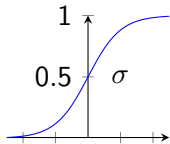
Vanishing Gradient Problem

Four neuron network, 1D input, 1D output



$$z_i = w_i a_{i-1} + b_i \quad a_i = \sigma(z_i)$$

$$\begin{aligned} \frac{\partial C}{\partial w_1} &= \frac{\partial C}{\partial z_4} \frac{\partial z_4}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \frac{\partial z_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1} \\ &= \sigma'(z_4) w_4 \cdot \sigma'(z_3) w_3 \cdot \sigma'(z_2) w_2 \cdot \sigma'(z_1) a_0 \\ &= \sigma'(z_4) \sigma'(z_3) \sigma'(z_2) \sigma'(z_1) \cdot w_4 w_3 w_2 a_0 \\ &= \underbrace{\sigma'(z_4) \sigma'(z_3) \sigma'(z_2) \sigma'(z_1)}_{\leq \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{4}} \cdot w_4 w_3 w_2 a_0 \end{aligned}$$

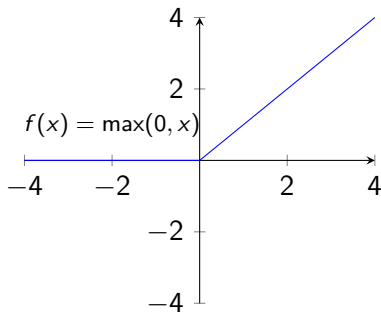


Vanishing Gradient Problem

- ▶ Derivative of sigmoidal activation functions < 1
- ▶ Exponential decay of gradient magnitude

Desirable: Activation function with derivative = 1 but non-linear
(> 1 = exploding gradient)

Solution: Rectified Linear Unit (ReLU) [Glorot & Bengio 2010]



How to Get Along With Limited Training Data?

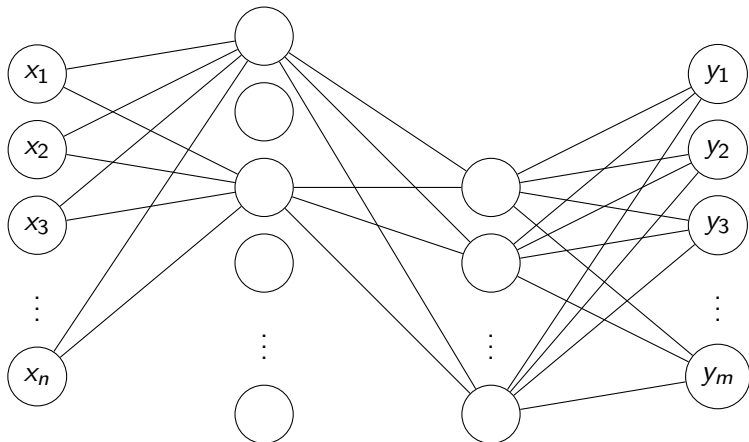
Problem: CNNs easily contain billions of parameters (weights)!
⇒ Could easily learn training samples “by heart”.

Solution: Apply *Regularization* during training

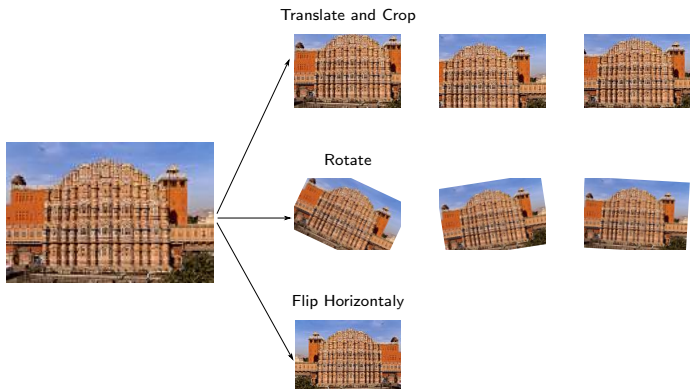
Fundamental Techniques:

- ▶ *Convolutional layers*
- ▶ *Dropout*
Randomly set outputs of neurons to zero
(usually 50% of fully connected layers)
- ▶ *Data Augmentation:*
Generate new, slightly different training samples from
existing ones by certain transforms
(e.g. slight translations, rotations, ...)

Dropout

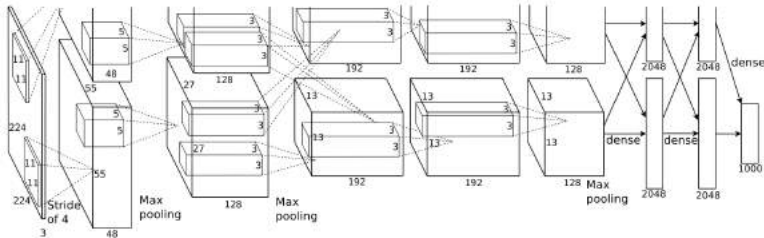


Data Augmentation



Note: Usually different augmentation techniques are mixed to create a single augmented image

Well-known Deep Learning Architectures: AlexNet

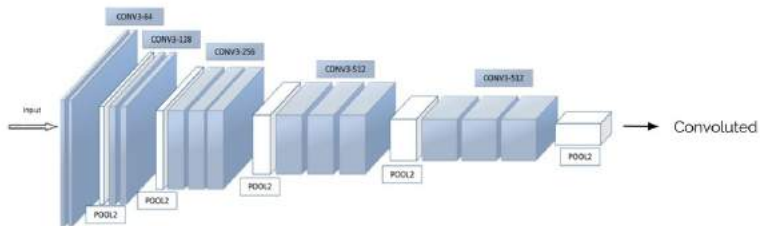


(Source: [Krizhevsky et al., 2012])

- ▶ CNN which kicked off the current Deep Learning hype
- ▶ Architecture similar to LeNet but more layers/parameters
- ▶ Trained on two graphic cards for over a week on ImageNet

A. Krizhevsky, I. Sutskever, G. E. Hinton: *ImageNet Classification with Deep Convolutional Neural Networks*, Neural Information Processing Systems, pp. 1097–1105, 2012.

Well-known Deep Learning Architectures: VGGNet



(Source: <http://html.scrip.org/>)

- ▶ First CNN to use only 3×3 convolutions (standard for current CNNs)
- ▶ Low number of filters in the early layers, high number of filters in the later layers
- ▶ Anytime pooling is applied, the number of filters is doubled

K. Simonyan, A. Zisserman: [Very Deep Convolutional Networks for Large-Scale Image Recognition](#), arXiv, 2014.

Summary Deep Learning

Some notes on Deep Learning:

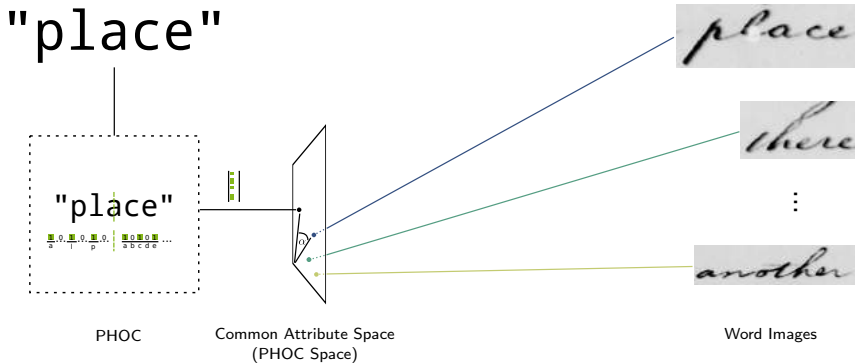
- ▶ Deep Neural Networks give state of the art results on a number of benchmarks across different tasks and domains
- ▶ Only empirical evidence that Deep Nets are better than shallower ones
- ▶ Suitable if hierarchy is present in your data
- ▶ Regularize where suitable/possible

(For a recent introduction to Deep Learning see [LeCun *et. al*, 2015])

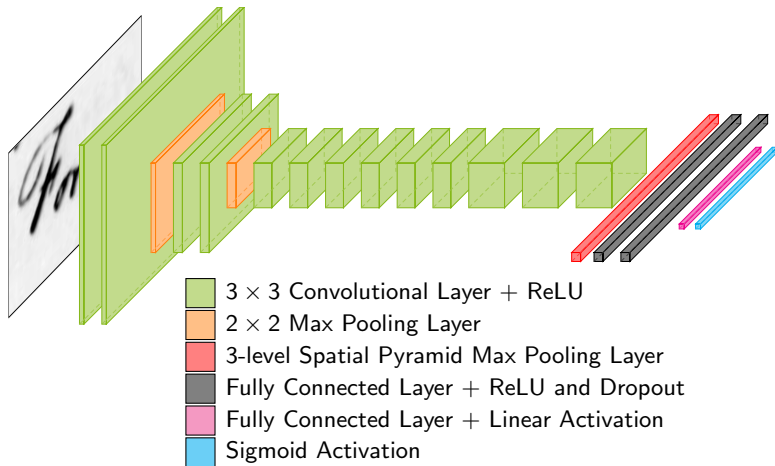
Overview

- ▶ Introduction
- ▶ Bag-of-Features: Fundamentals
- ▶ Learning Document Image Representations
- ▶ Learning Word Spotting Models
- ▶ Discussion
- ▶ Deep Learning Fundamentals
- ▶ **Deep Learning for Word Spotting**
- ▶ Summary

Reminder: General Framework



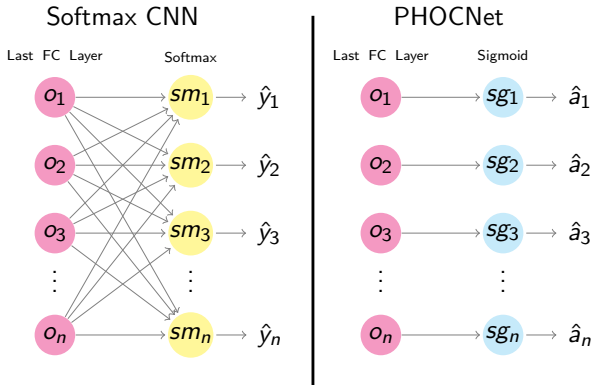
PHOCNet



S. Sudholt, G. A. Fink: [PHOCNet: A Deep Convolutional Neural Network for Word Spotting in Handwritten Documents](#), Proc. Int. Conf. on Frontiers in Handwriting Recognition, Shenzhen, China, 2016.

Softmax CNN vs. PHOCNet

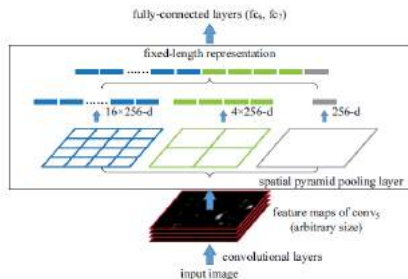
- ▶ In order to classify attributes, replace softmax with a sigmoid activation
- ▶ Each output neuron predicts one attribute
- ▶ Similar to *Logistic Regression*



Spatial Pyramid Pooling Layer

- ▶ Convolutional layers can already deal with arbitrary image sizes
- ▶ Only MLP part has a problem with changing image sizes

Solution: apply spatial pyramid pooling concept to the last convolutional output to generate fixed-size representation



K. He, X. Zhang, S. Ren, J. Sun: [Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition](#), Proc. European Conference on Computer Vision, pp. 346–361, 2014.

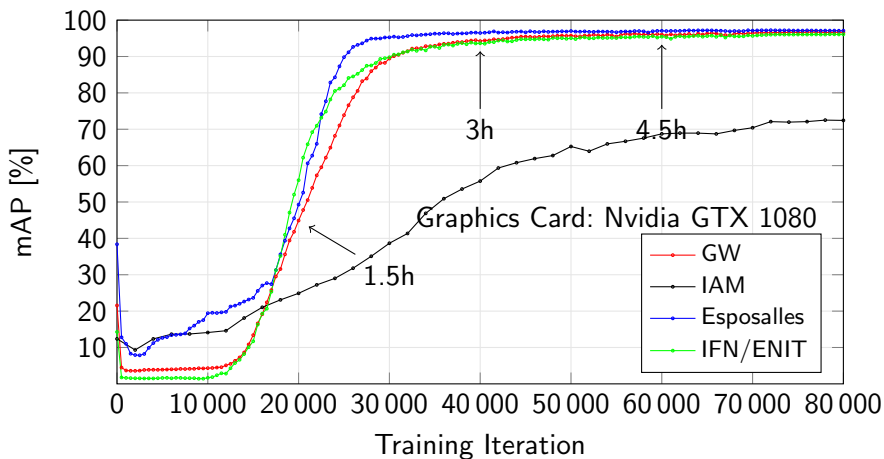
Experimental Evaluation

Segmentation-based Word Spotting Performance in mAP [%]

Method	GW		IAM		Esposalles		IFN/ENIT	
	QbE	QbS	QbE	QbS	QbE	QbS	QbE	QbS
PHOCNet	97.85	97.65	83.38	92.59	96.93	94.33	95.63	93.51
<i>Deep Feat. Emb.</i> [7]	94.41	92.84	84.24	91.58	-	-	-	-
Attribute SVM [2]	93.04	91.29	55.73	73.72	-	-	-	-
Finetuned CNN [23]	-	-	46.53	-	-	-	-	-
LSA Embedding [1]	-	56.54	-	-	-	-	-	-
BLSTM [3]	-	84.00	-	78.00	-	-	-	-
SC-HMM [16]	53.10	-	-	-	-	-	41.60	-
<i>Triplet-CNN</i> [27]	98.00	93.69	81.58	89.49	-	-	-	-

Experimental Evaluation II

mAP over the course of training

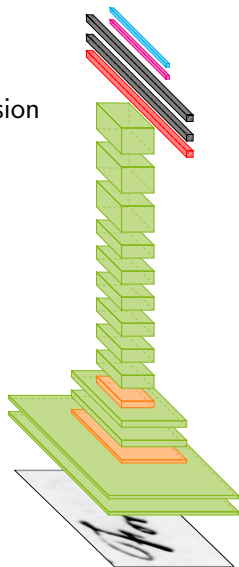


Overview

- ▶ Introduction
- ▶ Bag-of-Features: Fundamentals
- ▶ Learning Document Image Representations
- ▶ Learning Word Spotting Models
- ▶ Discussion
- ▶ Deep Learning Fundamentals
- ▶ Deep Learning for Word Spotting
- ▶ **Summary**

Summary

- ▶ Deep Learning is a hot topic in Computer Vision *and* Document Image Analysis
 - ⇒ *You should know about it!*
- ▶ CNNs are the most popular deep networks
 - ⇒ *Especially suitable for (document) images!*
- ▶ CNNs *can* work on limited data sets!
 - ⇒ *Always remember the PHOCNet :-)*
- ▶ Toolboxes (e.g. Caffe) make it easy to apply Deep Learning to ones own problems
 - ⇒ *Beware of blindly using them!*



References I

- [1] David Aldavert, Marçal Rusinol, Ricardo Toledo, and Josep Lladós.
Integrating visual and textual cues for query-by-string word spotting.
In International Conference on Document Analysis and Recognition, pages 511–515, 2013.
- [2] Jon Almazán, Albert Gordo, Alicia Fornés, and Ernest Valveny.
Word spotting and recognition with embedded attributes.
IEEE Trans. on Pattern Analysis and Machine Intelligence, 36(12):2552–2566, 2014.

References II

- [3] Volkmar Frinken, Andreas Fischer, R. Manmatha, and Horst Bunke.
A novel word spotting method based on recurrent neural networks.
IEEE Transactions on Pattern Analysis and Machine Intelligence, 34:211–224, 2012.

- [4] Xavier Glorot and Yoshua Bengio.
Understanding the difficulty of training deep feedforward neural networks.
AISTATS, 9:249–256, 2010.

References III

- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.
Spatial pyramid pooling in deep convolutional networks for visual recognition.
European Conference on Computer Vision, pages 346–361, 2014.
- [6] Kurt Hornik, Maxwell Stinchcombe, and Halbert White.
Multilayer feedforward networks are universal approximators.
Neural Networks, 2(5):359–366, 1989.
- [7] Praveen Krishnan, Kartik Dutta, and C.V. Jawahar.
Deep feature embedding for accurate recognition and retrieval of handwritten text.
In International Conference on Frontiers in Handwriting Recognition, pages 289–294, 2016.

References IV

- [8] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Neural Information Processing Systems*, pages 1097–1105, 2012.
- [9] Svetlana Lazebnik, Cordelia Schmid, and Jean Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *Proc. IEEE Comp. Soc. Conf. on Computer Vision and Pattern Recognition*, volume 2, pages 2169–2178, 2006.
- [10] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521:436–444, 2015.

References V

- [11] Yann LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel.
Handwritten digit recognition with a back-propagation network.
Neural Information Processing Systems, pages 396–404, 1990.
- [12] David Lowe.
Distinctive image features from scale-invariant keypoints.
Int. J. of Computer Vision, 60(2):91–110, 2004.
- [13] R. Manmatha, Chengfeng Han, E. M. Riseman, and W. B. Croft.
Indexing handwriting using word matching.
In *Proc. of the First ACM Int. Conf. on Digital Libraries*, DL '96,
pages 151–159, New York, NY, USA, 1996. ACM.

References VI

- [14] Stephen O'Hara and Bruce A. Draper.
Introduction to the bag of features paradigm for image classification and retrieval.
Computing Research Repository, arXiv:1101.3354v1, 2011.
- [15] Toni M. Rath and R. Manmatha.
Word image matching using dynamic time warping.
In *Computer Vision and Pattern Recognition, 2003. Proceedings. 2003 IEEE Computer Society Conference on*, volume 2, pages II-521-II-527 vol.2, June 2003.

References VII

- [16] José A. Rodríguez-Serrano and Florent Perronnin.
A model-based sequence similarity with application to
handwritten word spotting.
IEEE Transactions on Pattern Analysis and Machine Intelligence,
34(11):2108–2120, 2012.
- [17] F. Rosenblatt.
The perceptron: A probabilistic model for information storage
and organization in the brain.
Psychological Review, 65(6):386–408, 1958.

References VIII

- [18] Leonard Rothacker and Gernot A. Fink.
Segmentation-free query-by-string word spotting with
bag-of-features HMMs.
In Proc. Int. Conf. on Document Analysis and Recognition,
Nancy, France, 2015.

- [19] Leonard Rothacker, Marcal Rusinol, and Gernot A. Fink.
Bag-of-features HMMs for segmentation-free word spotting in
handwritten documents.
In Proc. Int. Conf. on Document Analysis and Recognition,
Washington DC, USA, 2013.

References IX

- [20] Leonard Rothacker, Marçal Rusinol, Josep Lladós, and Gernot A. Fink.
A Two-Stage Approach to Segmentation-Free Query-by-Example Word Spotting.
manuscript cultures, 1(7):47–57, 2014.
- [21] Leonard Rothacker, Szilard Vajda, and Gernot A. Fink.
Bag-of-features representations for offline handwriting recognition applied to Arabic script.
In Proc. Int. Conf. on Frontiers in Handwriting Recognition, Bari, Italy, 2012.

References X

- [22] M. Rusiñol, D. Aldavert, R. Toledo, and J. Lladós.
Browsing heterogeneous document collections by a
segmentation-free word spotting method.
In Proc. Int. Conf. on Document Analysis and Recognition, pages
63–67, Beijing, China, 2011.
- [23] Arjun Sharma and Sankar K. Pramod.
Adapting off-the-shelf CNNs for word spotting & recognition.
In International Conference on Document Image Analysis, pages
986–990, 2015.
- [24] R. Shekhar and C.V. Jawahar.
Word image retrieval using bag of visual words.
*In Document Analysis Systems (DAS), 2012 10th IAPR
International Workshop on*, pages 297–301, March 2012.

References XI

- [25] Karen Simonyan and Andrew Zisserman.
Very deep convolutional networks for large-scale image recognition.
arXiv, pages 1–13, 2014.
- [26] Sebastian Sudholt and Gernot A. Fink.
PHOCNet: A deep convolutional neural network for word spotting in handwritten documents.
In Proc. Int. Conf. on Frontiers in Handwriting Recognition, Shenzhen, China, 2016.

References XII

- [27] Tomas Wilkinson and Anders Brun.
Semantic and verbatim word spotting using deep neural networks.
In International Conference on Frontiers in Handwriting Recognition, pages 307–312, 2016.