# Universität Potsdam

## Institut für Informatik
## Lehrstuhl Maschinelles Lernen

# Sprachtechnologie
# Language Technology

Tobias Scheffer

Uwe Dick

# Organization

- Lecture and exercise
- Lecture: Monday, 12:15-13:45, 03.04.0.02
- Exercise: Monday, 14:15-15:45, 3.04.0.02
- Homework is posted on our website each Tuesday and is due the following Monday.

# Exam

- Successful completion of 70% of the homework exercises.
- 60 minutes of written exam, immediately followed by 15 mins of oral discussion.

# Organization

- Website
  - Contains slides and in some cases online lecture videos.

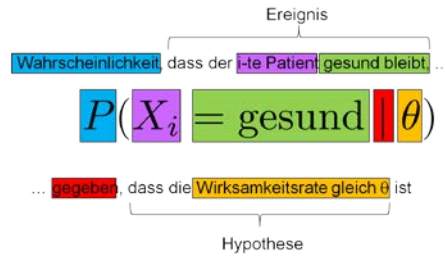# Literature

- Statistical natural language processing:
  - Manning & Schütze: „Foundations of Statistical Natural Language Processing." MIT Press
- Speech recognition
  - „The HTK Book", im Internet verfügbar.
    - Speech Recognition Toolkit
    - http://htk.eng.cam.ac.uk/docs/docs.shtml
  - Huang, Acero und Hon: „Spoken Language Processing". Prentice Hall.
- Information Retrieval:
  - Manning, Raghavan, Schütze: „Introduction to Information Retrieval". Cambridge University Press.

# Content

- Review of base technologies
  - Language models
  - Hiden Markov models, PCFG
  - Deep neural networks
- Language-processing tasks
  - Speech recognition
  - Translation, natural description of images
  - Parsing, information extraction
  - Text classification, clustering
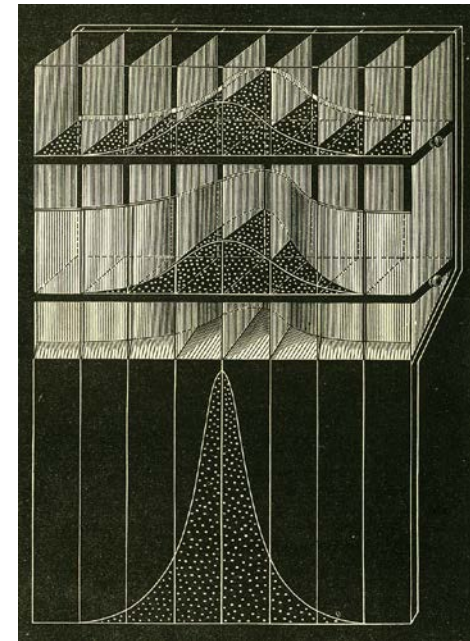- Information retrieval
  - Indexing and search

# Mathematical Foundations



$$P(X_i = \text{gesund} \mid \theta)$$

$$P(Y|X) = P(X|Y)\frac{P(Y)}{P(X)}$$

$$P(X_{neu}|X_1, \ldots, X_n) = \int_\theta P(X_{neu}|\theta, X_1, \ldots, X_n)P(\theta|X_1, \ldots, X_n)d\theta$$
$$= \int_\theta P(X_{neu}|\theta) \, P(\theta|X_1, \ldots, X_n) \, d\theta$$

# Statistical Language Models

- Elementary tool for
  - Speech recognition
  - Spell checking
  - Auto-complete, machine translation, ….

- Quantifies probability of a sequence of terms in a language, with respect to a corpus

$$P(w_1,...,w_T) = P(w_1)P(w_2 \mid w_1)...P(w_T \mid w_{T-1},...,w_1)$$

$$= P(w_1)P(w_2 \mid w_1)...P(w_T \mid w_{T-1}, w_{T-N+1})$$

$$= \prod_{i=1}^{N-1} P(w_i \mid w_{i-1},...,w_1) \prod_{i=N}^{T} P(w_i \mid w_{i-1},..w_{i-N+1})$$

# Statistical Language Models

- **Grammar, acceptor, parser**
  - Defines set of sentences of a language.
  - Defines hard boundaries.
  - Not a suitable mechanism for natural language.
  - Natural language has no hard boundaries, almost anything is possible.
- **Statistical language model**
  - Quantifies the probability of a sentence.
- **Statistical parser**
  - Infers the most likely syntactic interpretations of a sentence

# Statistical Language Models

- **N-Gram models**
  - ◆ Sequence model with n-th order Markov assumption
- **PCFG: probabilistic context-free grammar**

NP

Die          N'

$$\alpha_j(p,q) = P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} \mid G)$$

„Outside"

große          N'

braune          N'

$$\beta_j(p,q) = P(w_{pq} \mid N_{pq}^j, G)$$

„Inside"          Kiste

# Hidden-Markov-Modell

- Probabilistic model for sequences.
- Used, for instance, in
  - speech recognition,
  - part-of-speech tagging.

# Deep Neural Networks

- Computational models of neural information processing.

- Good at learning abstract feature representations of complex input objects

# Deep Neural Networks

- State of the art for
  - Speech recognition
  - Translation
  - Natural description of images
  - …



Input signals

Weighted input signals are aggregated

Axon: output signal

Synaptic weights: strengthened and weakened by learning processes

# Speech Recognition

- **Statistical speech recognition**
  - ◆ Acoustic model +
  - ◆ Language model.
- **Neural networks for speech recognition**
- **Text to speech**

# Machine Translation

# Text Generation, Description of Images

# POS-Tagging, Named Entity Recognition

- Identifying parts of speech
- Identifying proper names of persons, organizations, places, times, dates, genes, molecules, …

# Parsing

- Finding the most likely syntactical structure of sentences.
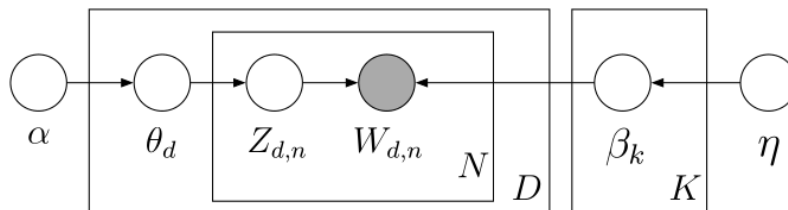
# Text Classification, Information Extraktion

# Clustering and Topic Models

- Clustering: grouping related documents.
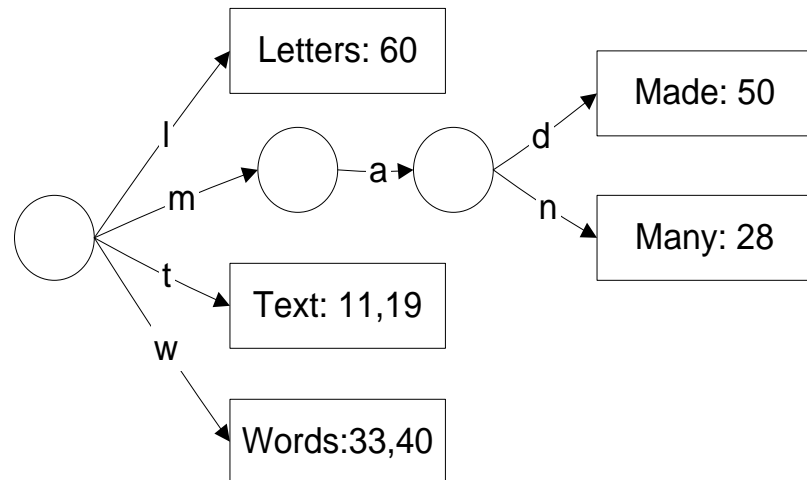- Topic modeling: groupung related terms in documents.

# Indexing and Search

- Fast search in large text corpora.

| 1 | 6 9 11 | 17 19 24 | 28 | 33 | 40 | 46 50 | 55 | 60 |

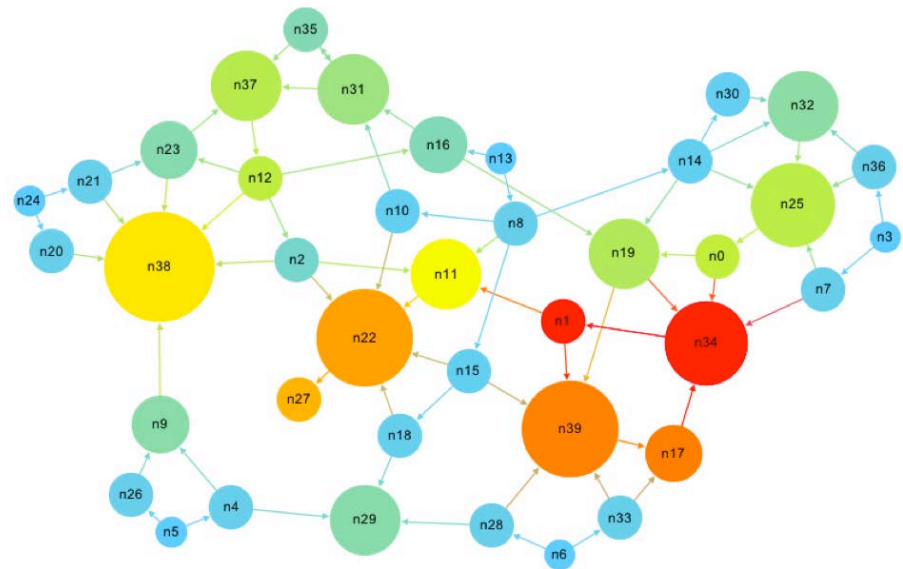This is a text. A text has many words. Words are made from letters.

| Terme | Vorkommen |
|---|---|
| Letters | 60 |
| Made | 50 |
| Many | 28 |
| Text | 11, 19 |
| words | 33, 40 |



21

# Web Search

- Crawling: visit which URL when? Challenges:
  - Infinite URLs, dynamic page content.
  - Update frequency and schedules.
  - Identical pages with multiple URLs.
  - Link spam.
- Relevance ranking, link analysis.

# Questions?

- No exercise sheet this week.
- No exercise this week.