# Statistical Natural Language Processing

Prof. Dr. Dietrich Klakow

# Lecture

- Lecture:
  - Friday 8:30-10:00
  - Location: MS Teams
  - Contact:
    - D. Klakow: dietrich.klakow@lsv.uni-saarland.de

# Exercises

- Exercises:
  - First exercise sheet will be issued today
  - Three groups:
    - Vilém Zouhar
      - vzouhar@lsv.uni-saarland.de
      - Thursday 16:00-17:30
    - Awantee Deshpande
      - adeshpande@lsv.uni-saarland.de
      - Tuesday 14:15-15:45
    - Julius Steuer
      - jsteuer@coli.uni-saarland.de
      - Tuesday 12:15-13:45

# Exercises

Submissions: Groups of 2 mandatory

Requirements:

- Threshold 70% (out of mandatory assignments)

- 10pts per assignment (around ~10) + bonus exercises (~2pts)

- 1pt bonus for active participation (showing up and **talking/contributing**)

- 5pts bonus for presenting the solution for an entire exercise sheet (at most once). Fractions for presenting parts of a sheet.

# Exercises

Submissions:

- Distribution of assignments using notebooks

- Suggestions: import your implementations rather than put it in the cells.

- Theoretical solutions also in the notebook (you can write LaTeX in cells)

# Teams + Piazza

- Lectures and tutorials will be on MS Teams
- Forum: Piazza

# Mailing List



If you haven't registered do so by today 11:00

# Exam

- Friday, July 23rd, 8:00-10:00.
- GHH E22, Foyer GHH and HS002 in E13.
- Assignment to locations will follow in the week before the exam once July corona rules are clear

# Literature

**Foundations of Statistical Natural Language Processing**

by <span style="color:red">Christopher D. Manning</span>, <span style="color:red">Hinrich Schütze</span>

**Publisher:** The MIT Press; 1st edition (June 18, 1999)

**ISBN:** 0262133601

**List Price:** $77.00



FOUNDATIONS OF STATISTICAL NATURAL LANGUAGE PROCESSING

CHRISTOPHER D. MANNING AND HINRICH SCHÜTZE

# Rules of the Game

- In case you don´t understand something:
    1. Ask!!!
    2. Ask!!!
    3. Ask!!!

# 1. Introduction: Overview of the Topics

# Use Zipf´s-Law in Language Modeling



Das Zipfsche Gesetz

# Chapter 2: Natural Language as a Sequence of Symbols

- Zipf's law

- Revision of basics of probability theory

# Guess the next word



President Bill ???

# Chapter 3: Basics of Language Modeling

- Language models for speech recognition
- Perplexity

# Coding a language efficiently

# Chapter 4: Entropy

- The Shannon game
- Text compression

# Chapter 5: Backing-Off Language Models

- Smoothing techniques

# Text Categorization



Speech Recognition

Information Retrieval

Computer Linguistics

Everything else

# Spam-Mail Classification

V / a g r a   $ 3 , 3 1

A m b / e n

M e r / d i a

C / a l i s   $ 3 , 7 5

V a l / u m   $ 1 , 2 1

X & n a x

S o m &

http://www.Chanatanxte.scriptmania.com/

# Chapter 6. Text Classification

- Variants of Task
- Algorithms
  - Nearest Neighbor Classifier
  - Maximum Entropy Models
  - Decision Trees
  - Neural Networks
  - Unsupervised Clustering

# Translation of the word „band" into German (output from LEO's)

- **band** das Band
- **band** die Band - *Musikgruppe*
- **band** [tech.] das Band
- **band** die Bandbreite
- **band** [chem.] die Bande - *im Spektrum* **band** das Beffchen **band** der Bereich **band** der Bund **band** der Frequenzbereich **band** die Gruppe **band** der Gurt **band** die Kapelle **band** die Leiste **band** die Musikkapelle **band** das Orchester **band** die Schar **band** die Schnur **band** [mus.] der Spielmannszug **band** der Streifen **band** die Truppe narrowband *also:* narrow-**band** *adj.* engbandig narrowband *also:* narrow-**band** *adj.* schmalbandig sideband *also:* side **band** [elec.] [telecom.] das Seitenband **Verben und Verbzusammensetzungen** to **band** together sich verbinden to **band** together sich vereinigen to **band** together sich zusammenrotten to **band** together sich zusammentun to **band** together zu einer Gruppe vereinigen to beat the **band** nie da gewesen sein to cross-**band** [tech.] absperren *[Holzverarbeitung]* **Zusammengesetzte Einträge** abrasive **band** - *cloth* [tech.] das Bandschleifleinen abrasive **band** - *paper* [tech.] das Bandschleifpapier adhesive **band** [tech.] das Klischeeklebeband attenuating **band** [aviat.] der Dämpfungsbereich audio **band** [phys.] der Hörbereich **band** aerial die Bandantenne **band**-aid das Heftpflaster **band**-aid [Amer.] [med.] das Pflaster **band**-aid [Amer.] [med.] das Wundpflaster **band** box die Hutschachtel **band** ceramics die Bandkeramik **band** collar der Stehkragen **band**-conveyor das Fließband **band** conveyor [tech.] der Gurtförderer **band**-conveyor das Transportband **band** edge die Bandkante **band** emission [autom.] die Bandemission **band** emission [autom.] die Bandenemission **band** gap [phys.] die Bandlücke **band** gate [tech.] der Bandausschnitt - *Spritzgusswerkzeug [Kunststoffe]* **band** grinder [tech.] die Bandschleifmaschine **band** matrix [math.] die Bandmatrix **band** of barrel das Fassband **band** of barrel der Fassreifen **band** of radiation [phys.] der Strahlungsbereich **band** of robbers die Räuberbande **band** overlap [tech.] die Bandüberlappung **band** printer [print.] der Banddrucker **band** radiation [autom.] die Bandenstrahlung **band** resaw [tech.] die Trennbandsäge **band** saw [tech.] die Bandsäge **band**-saw die Bandsäge **band** spectrum [tech.] das Bandenspektrum **band**-spread die Bandspreizung **band**-stand der Musikpavillon **band** structure [phys.] die Bandstruktur **band**-switch der Bereichsschalter **band**-switch der Bereichsumschalter **band** width die Bandbreite base **band** [tech.] das Basisband brake **band** [tech.] das Bremsband brass **band** [mus.] die Blaskapelle brass **band** [mus.] die Blechmusik brass **band** [mus.] der Spielmannszug broad **band** [tech.] das Breitband carrier **band** [tech.] das Trägerfrequenzband clay **band** [geol.] das Salband clincher **band** [autom.] das Wulstband *[Reifen]* conveyer **band** das Förderband cover **band** [tech.] das Deckband currency **band** [bank.] die Währungsbandbreite dance **band** die Tanzkapelle dead **band** [metr.] die Totzone edge **band** [tech.] der Umleimer *[Tischlerei]* elastic **band** [tech.] das Gummiband elastic **band** der Gummistrumpf error **band** der Zufallsstreubereich filter **band** [tech.] das Siebband flexible **band** die Randzeit - *Arbeitszeit* glassy **band** [tech.] glasiger Streifen guard **band** [elec.] der Rasen - *Abstand zwischen den Schrägspuren, den Videospuren, der benutzt wird, um eine gegenseitige Beeinflussung der Spuren zu vermeiden.* guard **band** [elec.] der Schutzabstand - *Abstand zwischen den Schrägspuren, den Videospuren, der benutzt wird, um eine gegenseitige Beeinflussung der Spuren zu vermeiden.* guard **band** [elec.] [telecom.] der Schutzbereich - *zwischen zwei Kanälen zur Vermeidung von Interferenzen.* guard **band** [elec.] [telecom.] der Schutzbereicht guard **band** [elec.] [telecom.] das Sicherheitsband - *zwischen zwei Kanälen zur Vermeidung von Interferenzen.* guard **band** [elec.] [telecom.] das Sicherheitsfrequenzband - *zwischen zwei Kanälen zur Vermeidung von Interferenzen.* guide **band** das Führungsband hair-**band** das Haarband heating **band** [tech.] das Heizband hinge **band** [tech.] das Gelenkband mehr >>

# Chapter 7: Word Sense Disambiguation

- Dictionary-Based Disambiguation
- Thesaurus based methods
- Bayes Classifier

# Example for Part-Of-Speech Tagging

Xinhua News Agency , Guangzhou , March 16 ( Reporter Chen Ji ) The latest statistics show that from January through February this year , the export of high-tech products in Guangdong Province reached 3.76 billion US dollars , up 34.8% over the same period last year and accounted for 25.5% of the total export in the province .

# Example for Part-Of-Speech Tagging

Xinhua/NNP News/NNP Agency/NNP ,/, Guangzhou/NNP ,/, March/NNP 16/CD (/( Reporter/NNP Chen/NNP Ji/NNP )/SYM The/DT latest/JJS statistics/NNS show/VBP that/IN from/IN January/NNP through/IN February/NNP this/DT year/NN ,/, the/DT export/NN of/IN high-tech/JJ products/NNS in/IN Guangdong/NNP Province/NNP reached/VBD 3.76/CD billion/CD US/PRP dollars/NNS ,/, up/IN 34.8%/CD over/IN the/DT same/JJ period/NN last/JJ year/NN and/CC accounted/VBD for/IN 25.5%/CD of/IN the/DT total/JJ export/NN in/IN the/DT province/NN ./.

# Chapter 8: Part-Of-Speech Tagging

- Hidden Markov Model
- Rule based: The Brill tagger

# Chapter 9. Named Entity Tagging

Task:

Identify names of people, organizations, locations … in text

- President <ENAMEX id="9" type="PERSON">Richard Nixon</ENAMEX> in <ENAMEX id="10" type="LOCATION">Moscow.</ENAMEX>
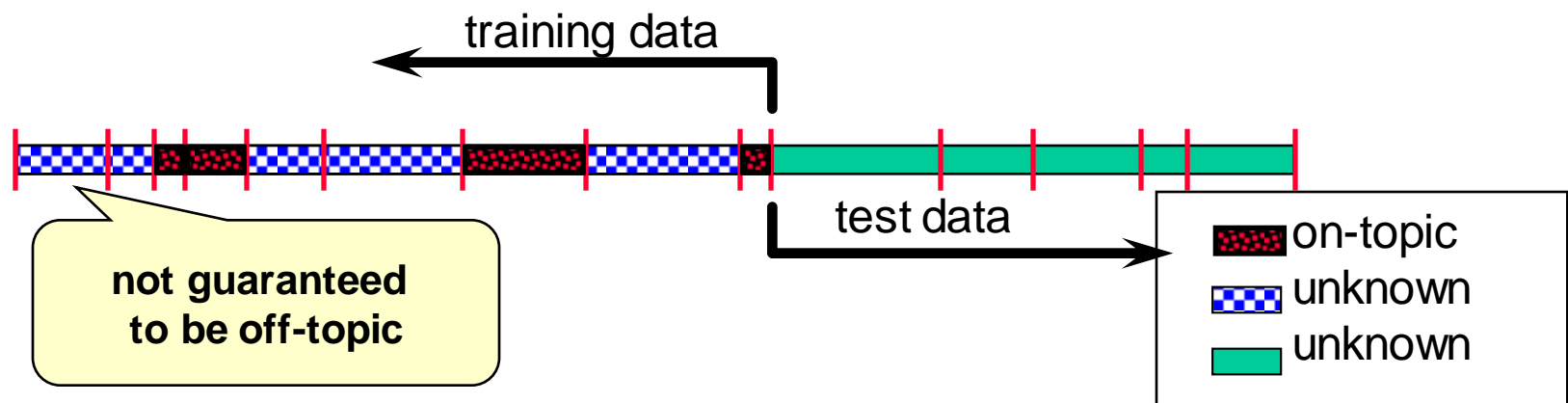
28

# Chapter 10: Information Retrieval

- Evaluation
- Processing the query
- Vector space model
- Term weighting
- Distance metrics
- Models for term distribution
- Probabilistic IR
- Singular value decomposition
- Language models

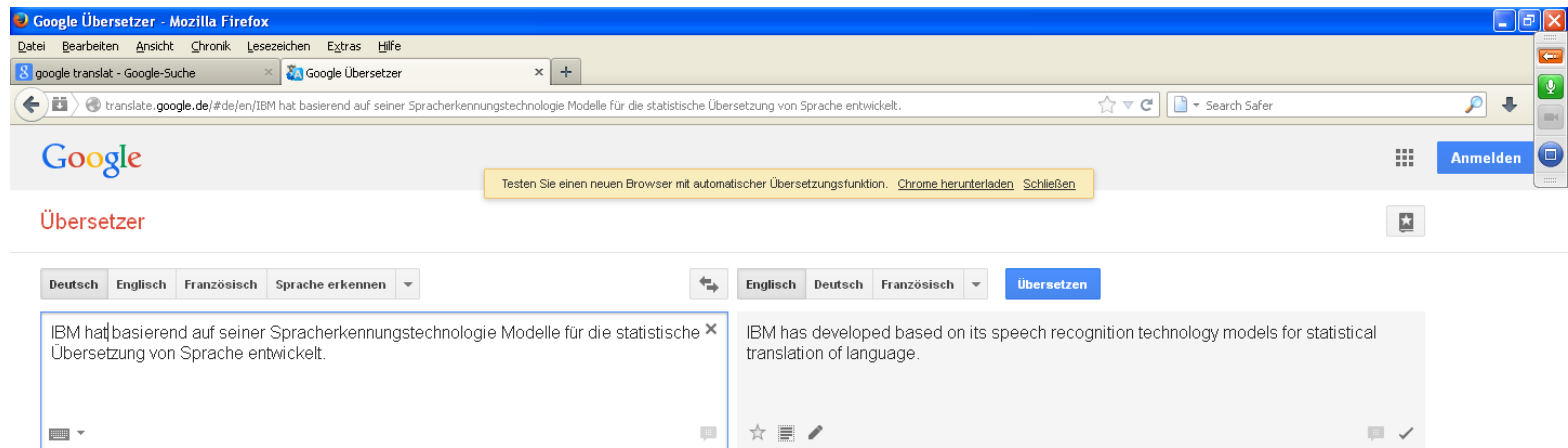# Chapter 11. Topic Detection and Tracking (if time permits)

*To detect stories that discuss the target topic, in multiple source streams.*

- Find all the stories that discuss a given target topic
  - *Training:* Given $N_t$ sample stories that discuss a given target topic,
  - *Test:* Find all subsequent stories that discuss the target topic.

training data

not guaranteed to be off-topic

test data

on-topic
unknown
unknown

# Chapter 12: Statistical Machine Translation

- Machine translation as a sequence labeling problem
- IBM models 1-4

# Summary Chapter 1

- Organization of the lecture
- Overview of the topics

- Are those the topics you are expecting?