# Information Retrieval & Data Mining

Universität des Saarlandes, Saarbrücken

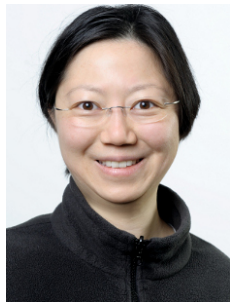Winter Semester 2013/14

# The Course

## Lecturers



Klaus Berberich
kberberi@mpi-inf.mpg.de



Pauli Miettinen
pmiettin@mpi-inf.mpg.de
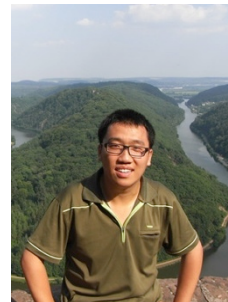
## Teaching Assistants



Amy Siu
sui@mpi-inf.mpg.de



Arunav Mishra
amishra@mpi-inf.mpg.de



Erdal Kuzey
ekuzey@mpi-inf.mpg.de



Kai Hui
khui@mpi-inf.mpg.de



Kaustubh Beedkar
kbeedkar@mpi-inf.mpg.de



Sourav Dutta
sdutta@mpi-inf.mpg.de

D5: Databases & Information Systems Group
Max Planck Institute for Informatics

# Organization

- **Lectures:**
  - **Tuesday 16-18** and **Thursday 14-16**
    in **Building E1.3**, **HS-002**

- **Office hours:**
  - **Tuesday 14-16**

- **Assignments/tutoring groups**
  - **Monday 12-14 / 14-16 / 16-18, R021, E1.4** (MPI-INF building)
  - **Friday 12-14 / 14-16, R021, E1.4** (MPI-INF building)

  Assignments given out in Thursday lecture, to be solved until next Thursday
  - First assignment sheet given out on **Thursday, Oct 17**
  - First meetings of tutoring groups on **Friday, Oct 25**

# Requirements for Obtaining 9 Credit Points

- **Pass 2 out of 3 written tests**
  Tentative dates: **Tue, Nov 12**; **Thu, Dec 12**; **Tue, Jan 28**
  (45-60 min each)

- **Pass the final written exam**

  Tentative date: **Tue, Feb 13** (120-180 min)

- Must **present solutions to 3 assignments**, more possible
  (**You must return your assignment sheet and have a correct
  solution  in order to present in the exercise groups.**)
  - **1 bonus point** possible in tutoring groups
  - **Up to 3 bonus points** possible in tests
  - Each bonus point earns one mark in letter grade
  (0.3 in numerical grade)

# Register for Tutoring Groups

http://bit.ly/irdm

- Register for one of the tutoring groups **until Oct 22**
- Check back frequently for updates & announcements

# Agenda

I.      Introduction

II.     Probability theory, statistics, linear algebra

II.     Ranking principles

III.    Link analysis

IV.     Indexing & searching

V.      Information extraction

VI.     Frequent itemsets & association rules

VII.    Unsupervised clustering

VIII.   (Semi-)supervised classification

IX.     Advanced topics in data mining

X.      Wrap-up & summary

**Information Retrieval**

**Data Mining**

# Literature (I)

- **Information Retrieval**

  - Christopher D. Manning, Prabhakar Raghavan, Hinrich Schütze.
    *Introduction to Information Retrieval*
    Cambridge University Press, 2008.
    Website: http://nlp.stanford.edu/IR-book/

  - R. Baeza-Yates, R. Ribeiro-Neto.
    *Modern Information Retrieval: The concepts and technology behind search.*
    Addison-Wesley, 2010.
    Website: http://www.mir2ed.org

  - W. Bruce Croft, Donald Metzler, Trevor Strohman.
    *Search Engines: Information Retrieval in Practice.*
    Addison-Wesley, 2009.
    Website: http://www.pearsonhighered.com/croft1epreview/

# Literature (II)

- **Data Mining**

  – Mohammed J. Zaki, Wagner Meira Jr.

  *Data Mining and Analysis: Fundamental Concepts and Algorithms*
  Manuscript (will be made available during the semester)

  – Pang-Ning Tan, Michael Steinbach, Vipin Kumar.

  *Introduction to Data Mining*

  Addison-Wesley, 2006.

  Website: http://www-users.cs.umn.edu/%7Ekumar/dmbook/index.php

# Literature (III)

- **Background & Further Reading**

  - Jiawei Han, Micheline Kamber, Jian Pei.
    *Data Mining - Concepts and Techniques*, 3rd ed., Morgan Kaufmann, 2011
    Website: http://www.cs.sfu.ca/~han/dmbook

  - Stefan Büttcher, Charles L. A. Clarke, Gordon V. Cormack.
    *Information Retrieval: Implementing and Evaluating Search Engines*,
    MIT Press, 2010

  - David B. Skillicorn.
    *Understanding complex datasets: data mining with matrix decomposition*,
    Chapman & Hall/CRC, 2007

  - Christopher M. Bishop.
    *Pattern Recognition and Machine Learning*, Springer, 2006

  - Larry Wasserman.
    *All of Statistics*, Springer, 2004
    Website: http://www.stat.cmu.edu/~larry/all-of-statistics/

# Quiz Time!

- Please answer the **20 quiz questions** during the rest of the lecture.

- The quiz is completely **anonymous**, but keep your id on the top-right corner. There will be a **prize for the 3 best answer** sheets.

# Chapter I:
# Introduction – Information Retrieval and Data Mining in a Nutshell

Information Retrieval & Data Mining

Universität des Saarlandes, Saarbrücken

Winter Semester 2013/14

# Chapter I: Information Retrieval and Data Mining in a Nutshell

- **1.1 Information Retrieval in a Nutshell**

  – Search & beyond

- **1.2 Data Mining in a Nutshell**

  – Real-world DM applications

*„We are drowning in information, and starved for knowledge.“*
*-- John Naisbitt*

# I.1 Information Retrieval in a Nutshell

- Web, intranet, digital libraries, desktop search
- Unstructured/semi-structured data

**crawl** → **extract & clean** → **index** → **match** → **rank** → **present**

handle dynamic pages, detect duplicates, detect spam

fast top-k queries, query logging, auto-completion

GUI, user guidance, personalization

strategies for crawl schedule and priority queue for crawl frontier

build and analyze web graph, index all tokens or word stems

scoring function over many data and context criteria

**Server farms** with **10 000's** (2002) – **100,000's** (2010) computers, distributed/replicated data in high-performance file system (**GFS**, **HDFS**,…), massive parallelism for query processing (**MapReduce**, **Hadoop**,…)

# Content Preprocessing



**Search Engines**
Politicians are worried that the Web is now dominated by search engine companies …

**Document**

politicians
worried
web
**...**

Extraction of **salient words**

politic
worry
web
**...**

Linguistic methods: **stemming**, **lemmas**

Statistically **weighted features** (terms)

*Thesaurus*

Synonyms, Sub-/Super-Concepts

politic
law
firm
worry
web
politic
web
search
…

**Bag of words**

# Vector Space Model for Relevance Ranking

**Ranking** by descending relevance

$\longleftarrow$
$\longleftarrow$
$\longleftarrow$

**Search engine**

**Query** $q \in [0,1]^{|F|}$
(feature vector)

Documents are **feature vectors**

**Similarity metric:**

$$sim(d_i, q) := \frac{\sum_{j=1}^{|F|} d_{ij}\, q_j}{\sqrt{\sum_{j=1}^{|F|} d_{ij}^2 \sum_{j=1}^{|F|} q_j^2}}$$

$with\ d_i \in [0,1]^{|F|}$

e.g., using:

$$d_{ij} := w_{ij} / \sqrt{\sum_k w_{ik}^2}$$

$$w_{ij} := \log\left(1 + \frac{freq(f_j, d_i)}{\max_k freq(f_k, d_i)}\right) \log \frac{\#docs}{\#docs\,with\,f_i}$$

**tf\*idf formula**

# Link Analysis for Authority Ranking

**Ranking** by descending **relevance & authority**

**Search engine**

**Query** $q \in [0,1]^{|F|}$
(feature vector)



+ **Consider in-degree and out-degree of web pages:**
   **Authority** $(d_i) :=$
      **Stationary visiting probability** $[d_i]$
      **in random walk on the Web (ergodic Markov Chain)**

+ **Reconciliation of relevance and authority by ad hoc weighting**

# Google's PageRank [Page and Brin 1998]

- **Ideas:** (i) Hyperlinks are endorsements
  (ii) Page is important if many important pages link to it

- **Random walk** on web graph $G(V, E)$ with random surfer that randomly follows outgoing link or jumps to another random page

$$P(v) = (1 - \varepsilon) \sum_{(u,v) \in E} \frac{P(u)}{out(u)} + \frac{\varepsilon}{|V|}$$

- **PageRank** $P(v)$ corresponds to the stationary visiting probability of state $v$ in an ergodic Markov chain

# Inverted Index

Vector space model suggests **term-document matrix**,
but data is sparse and queries are even very sparse
→ better use **inverted index** with terms as keys for B+ tree

**q: professor**
   **research**
   **xml**

B+ tree on terms

| professor | ••• | research | ••• | xml |

**index lists**
**with postings**
**(DocId, Score)**
**sorted by DocId**

| professor | research | xml |
|-----------|----------|-----|
| 17: 0.3   | 12: 0.5  | 11: 0.6 |
| 44: 0.4   | 14: 0.4  | 17: 0.1 |
| 52: 0.1   | 28: 0.1  | 28: 0.7 |
| 53: 0.8   | 44: 0.2  | ⋮ |
| 55: 0.6   | 51: 0.6  | |
| ⋮         | 52: 0.3  | |
|           | ⋮        | |

**Google:**
> 10 Mio. terms
> 20 Bio. docs
> 10 TB index

terms can be full words, word stems, word pairs, substrings, N-grams, etc.
(whatever "dictionary terms" we prefer for the application)

- index-list entries in **DocId order** for fast Boolean operations
- many techniques for excellent **compression** of index lists
- additional **position index** needed for phrases, proximity, etc.
  (or other pre-computed data structures)

# Evaluation of Search Result Quality

Ideal measure is "**satisfaction of user's information need**" heuristically approximated by benchmarking measures (on test corpora with query suite and relevance assessment by experts)

Capability to return **only** relevant documents:

$$\textbf{\textit{Precision}} = \frac{\text{\# relevant docs among top r}}{r}$$

typically for r = 10, 100, 1000

Capability to return **all** relevant documents:

$$\textbf{\textit{Recall}} = \frac{\text{\# relevant docs among top r}}{\text{\# relevant docs}}$$

typically for r = corpus size

**Typical quality**

**Ideal quality**

# Beyond Web Search…

- Find answers to **"knowledge queries"** and **natural language questions** (e.g., by scientists or journalists)
  - *Who was German chancellor when Angela Merkel was born?*
  - *How are Max Planck, Angela Merkel, and the Dalai Lama related?*
  - *Which politicians are also entrepreneurs?*
  - *What was the population of Munich in 1972?*
  - *…*

- Knowledge about **entities** (e.g., persons and locations), **classes**, **attributes**, **relationships** between them is required
  - focus on **structured data sources** (e.g., relational, XML, RDF)
  - perform **information extraction** on semi-structured & textual data

# Google Knowledge Graph



http://www.google.com

# Freebase



http://www.freebase.com

# YAGO

http://www.yago-knowledge.org

# DBpedia



**About: Dave Grohl**

An Entity of Type : musical artist, from Named Graph : http://live.dbpedia.org, within Data Space : live.dbpedia.org

David Eric "Dave" Grohl (born January 14, 1969) is an American rock musician, multi-instrumentalist, singer-songwriter and film director, who is the lead vocalist, guitarist, primary or main songwriter and founder of the band Foo Fighters. Prior to Foo Fighters, Grohl was the drummer for the grunge band Nirvana. He is also the drummer and co-founder of the rock supergroup Them Crooked Vultures.

| Property | Value |
| --- | --- |
| dbpedia-owl:abstract | • David Eric "Dave" Grohl (born January 14, 1969) is an American rock musician, multi-instrumentalist, singer-songwriter and film director, who is the lead vocalist, guitarist, primary or main songwriter and founder of the band Foo Fighters. Prior to Foo Fighters, Grohl was the drummer for the grunge band Nirvana. He is also the drummer and co-founder of the rock supergroup Them Crooked Vultures. Grohl has additionally written all the music and performed all the instruments for his short-lived side projects Late! and Probot, as well as being involved with Queens of the Stone Age numerous times throughout the past decade. He has performed session work (as a drummer) for a variety of musicians, including Garbage, Killing Joke, Nine Inch Nails, David Bowie, Paul McCartney, The Prodigy, Slash, Iggy Pop, Juliette Lewis, Tenacious D, Tom Petty and the Heartbreakers, Lemmy and Stevie Nicks. |
| dbpedia-owl:activeYearsStartYear | • 1984-01-01 00:00:00 (xsd:date)<br>• 1984-01-01 00:00:00 (xsd:date) |
| dbpedia-owl:alias | • Davy Grolton, Dale Nixon, Late! (pseudonym for his solo album Pocketwatch), and Dr. G (as Tenacious D's drummer) . |
| dbpedia-owl:associatedBand | • dbpedia:Paul_McCartney<br>• dbpedia:Stevie_Nicks<br>• dbpedia:The_Prodigy<br>• dbpedia:Trent_Reznor<br>• dbpedia:Tom_Petty_and_the_Heartbreakers<br>• dbpedia:Rick_Springfield<br>• dbpedia:Killing_Joke<br>• dbpedia:Probot<br>• dbpedia:Mondo_Generator<br>• dbpedia:Tenacious_D<br>• dbpedia:Foo_Fighters<br>• dbpedia:Queens_of_the_Stone_Age<br>• dbpedia:Them_Crooked_Vultures<br>• dbpedia:Dain_Bramage<br>• dbpedia:Nirvana_(band)<br>• dbpedia:Slash_(musician)<br>• dbpedia:Scream_(band) |
| dbpedia-owl:associatedMusicalArtist | • dbpedia:Paul_McCartney<br>• dbpedia:Stevie_Nicks<br>• dbpedia:The_Prodigy<br>• dbpedia:Trent_Reznor<br>• dbpedia:Tom_Petty_and_the_Heartbreakers<br>• dbpedia:Rick_Springfield<br>• dbpedia:Killing_Joke<br>• dbpedia:Probot<br>• dbpedia:Mondo_Generator<br>• dbpedia:Tenacious_D<br>• dbpedia:Foo_Fighters<br>• dbpedia:Queens_of_the_Stone_Age<br>• dbpedia:Them_Crooked_Vultures<br>• dbpedia:Dain_Bramage<br>• dbpedia:Nirvana_(band)<br>• dbpedia:Slash_(musician)<br>• dbpedia:Scream_(band) |
| dbpedia-owl:background | • solo_singer |
| dbpedia-owl:birthDate | • 1969-01-14 (xsd:date)<br>• 1969-01-14 (xsd:date) |
| dbpedia-owl:birthPlace | • dbpedia:United_States<br>• dbpedia:Ohio<br>• dbpedia:Warren,_Ohio<br>• dbpedia:Norrköping,_Sweden |
| dbpedia-owl:genre | • dbpedia:Hardcore_punk<br>• dbpedia:Alternative_rock<br>• dbpedia:Hard_rock<br>• dbpedia:Heavy_metal_music<br>• dbpedia:Post-grunge<br>• dbpedia:Grunge |

http://dbpedia.org

# The Linked Data Project



as of 2011:
- 295 sources
- 32 billion triples
- 504 million links

As of September 2010 (cc) (i) (s)

http://linkeddata.org

BBC Programmes · BBC Music · Chronicling America · Event-Media · Linked MDB · NSZL Catalog · PBAC · MARC Codes List · Semantic Crunch Base · Lotico · Revyu · SW Dog Food · (Tune)

BBC Wildlife Finder · Recht-spraak.nl · Tele-graphis · New York Times · URI Burner · flickr wrappr · OpenCalais · Good-win Family · RDF ohloh · BibBase · DBLP (L3S) · DBLP (RKB Explorer) · (R Exp

US · Taxon Concept · Geo Names · World Fact-book (FUB) · Freebase · DBpedia · iServe · VIVO UF · VIVO Indiana · VIVO Cornell · DBLP (FU Berlin) · IEEE · CiteS

NASA (Data Incubator) · Fishes of Texas · Geo Species · Uberblic · dbpedia lite · OS · data dcs · GESIS · Course-ware

Euro-stat (FUB) · Geo Linked Data (es) · UMBEL · lingvoj · YAGO · Daily Med · TCM Gene DIT · SIDER · Project Gutenberg (FUB) · ERA · STW · UN/LOCODE · Pub Chem

ked r Data e.sis) · riese · Open Cyc · Lexvo · totl.net · Linked CT · Uni Pathway · Drug Bank · Medi Care · Disea-some · STITCH · OBO · KE Dr · ChEBI · KEGG Cpd

UNIS · Twarql · WordNet (VUA) · Linked Open Numbers · UniRef · Taxo-nomy · UniProt · Pfam · PDB · Reactome · HGNC · CAS

Linked GeoData · WordNet (W3C) · Cornetto · UniParc · Affy-metrix · PRO-SITE · ProDom · PubMed · Gene Ontology · SGD · Chem2 Bio2RDF · Homo Gen

Airports · Product DB · UniSTS · Gen Bank · OMIM · InterPro · MGI · GeneID

# Jeopardy!

A big US city with two airports, one named after a World War II hero, and one named after a World War II battle field?



Chicago - Wikipedia, the free encyclopedia - Mozilla Firefox

O'Hare International Airport - Wikipedia, the free encyclopedia - Mozilla Firefox

Chicago Midway International Airport - Wikipedia, the free encyclopedia - Mozilla Firefox

File   Edit   View   History   Bookmarks   Tools   Help

W   http://en.wikipedia.org/wiki/Chicago_Midway_International_Airport

W Chicago Midway International Airport - Wikip...   +

Main page
Contents
Featured content
Current events
Random article
Donate to Wikipedia

Interaction
Help
About Wikipedia
Community portal
Recent changes
Contact Wikipedia

Toolbox

Print/export

Languages
Deutsch

## Chicago Midway International Airport

From Wikipedia, the free encyclopedia

*"MDW" redirects here. For other uses, see MDW (disambiguation).*

*For other uses, see Midway Airport (disambiguation).*

**Chicago Midway International Airport** (IATA: **MDW**, ICAO: **KMDW**, FAA LID: **MDW**), also known simply as **Midway Airport** or **Midway**, is an airport in Chicago, Illinois, United States, located on the city's southwest side, eight miles (13 km) from Chicago's Loop. The airport's current IATA code MDW has been used since 1949 when Chicago Municipal Airport was renamed Chicago Midway Airport,[3] although the airline schedule books continued to call it CHI until airline flights began at O'Hare. It is bordered by 55th Street, Cicero Avenue (terminal entrance), 63rd Street, and Central Avenue. The airport's northern half is within the Garfield Ridge community area, and the southern half is within the Clearing community area. The airport is managed by the Chicago Airport System, which also oversees operations at O'Hare International Airport and Gary/Chicago International Airport.[4] The airport is named after the Battle of Midway during World War II.

Midway is dominated by low-cost carrier Southwest Airlines. AirTran Airways and Delta Air Lines are the

Chicago M

Aerial view of
a.k.a. th

IATA: MD

# Deep-QA in NL

William Wilkinson's "An Account of the Principalities of Wallachia and Moldavia" inspired this author's most famous novel

This town is known as "Sin City" & its downtown is "Glitter Gulch"

As of 2010, this is the only former Yugoslav republic in the EU

99 cents got me a 4-pack of Ytterlig coasters from this Swedish chain

**question classification & decomposition** → **knowledge backends**

D. Ferrucci et al.: **Building Watson: An Overview of the DeepQA Project.** AI Magazine, 2010.
**www.ibm.com/innovation/us/watson/index.htm**

# IRDM Research Literature

Important **conferences** on IR and DM
(see DBLP bibliography for full detail, http://www.dblp.org)
SIGIR, WSDM, ECIR, CIKM, WWW, KDD, ICDM, ICML, ECML

Important **journals** on IR and DM
(see DBLP bibliography for full detail, http://www.dblp.org)
TOIS, TOW, InfRetr, JASIST, InternetMath, TKDD, TODS, VLDBJ

Performance **evaluation/benchmarking** initiatives:
• Text Retrieval Conference (TREC), http://trec.nist.gov
• Cross-Language Evaluation Forum (CLEF), http://www.clef-campaign.org
• Initiative for the Evaluation of XML Retrieval (INEX),
    http://www.inex.otago.ac.nz/
•  KDD Cup, http://www.kdnuggets.com/datasets/kddcup.html
                    & http://www.sigkdd.org/kddcup/index.php