![HPI Hasso Plattner Institut — IT Systems Engineering | Universität Potsdam]

# Search Engines
# Chapter 1 – Introduction

21.4.2009

Felix Naumann

# Brazil

- **Sam Lowry**: My name's Lowry. Sam Lowry. I've been told to report to Mr. Warren.
- **Porter - Information Retrieval**: Thirtieth floor, sir. You're expected.
- **Sam Lowry**: Um... don't you want to search me?
- **Porter - Information Retrieval**: No sir.
- **Sam Lowry**: Do you want to see my ID?
- **Porter - Information Retrieval**: No need, sir.
- **Sam Lowry**: But I could be anybody.
- **Porter - Information Retrieval**: No you couldn't sir. This is Information Retrieval.
- Sources
    - http://en.wikiquote.org/wiki/Brazil_(film)
    - http://www.youtube.com/watch?v=LFlFIG22Y9E&hl=de

# Anthropology Program at Kansas State University – Michael Wesch

- Information (r)evolution
    - http://www.youtube.com/watch?v=-4CV05HyAbM
    - http://ksuanth.weebly.com/wesch.html
- The machine is Us/ing us
    - http://www.youtube.com/watch?v=NLlGopyXT_g

# Overview

- Introduction to team

- Organization

- Information Retrieval & Search Engines

- Overview of semester

# Information systems team

HPI Hasso Plattner Institut

project **ViQTOR**

Katrin **Heinrich**

Prof. Felix **Naumann**

Jens **Bleiholder**

project **fusem**

Paul **Führing**

**DQ Annotation & Assessment**

**Data Fusion**

project **HumMer**

**Information Integration**

**Data Profiling & Cleaning**

Christoph Böhm

**Information Quality**

Armin **Roth**

**Peer Data Management Systems**

**Matching**

**Data Integration for Life Science Data Sources**

**ETL Management**

project **System P**

**Service-Oriented Systems**

project **Aladin**

**Ontologies, Profiling**

Mohammed **AbuJarour**

Frank **Kaufer**

Jana **Bauckmann**

Alexander **Albrecht**

**Data Profiling for Schema Management**

# Other courses in this semester

Lectures

- DBS I
- Search engines

Seminars

- Bachelor: Beauty is our Business
- Bachelor: Map/Reduce Algorithms on Hadoop
- Master: Linked Data Profiling
- Forschungsseminar

Bachelorproject

- ETL Management



Extending the Database Relational Model to Capture More Meaning

E. F. CODD
IBM Research Laboratory

# Overview

- Introduction to team
- Organization
- Information Retrieval & Search Engines
- Overview of semester

- Lectures
  - Tuesday 9:15 – 10:45
  - Thursdays 13:30 – 15:00
- Practical work
  - Selected dates – see webpage
- First lecture
  - 21.4.2009
- Last lecture
  - 23.7.2009
- Holidays
  - 21.5. Himmelfahrt

- Exam
  - Oral exam, 30 minutes
  - First week after lectures end
- 7 exercise courses
  - TAs: Alexander Albrecht & Mohammed AbuJarour
  - Practical work and presentations
  - Teams of two students
- Prerequisites
  - For participation
    - ◇ Basic knowledge in databases
  - For exam
    - ◇ Attendance of lectures
    - ◇ Active participation in exercise courses
    - ◇ Successful work on all practical assignments
      - – "Success" to be defined

# Feedback

- Evaluation at end of semester
- Q&A anytime!
    - During lecture
    - Directly after lecture
    - Consultation: Tuesday 15-16
    - Email: naumann@hpi.uni-potsdam.de
- Hints on improvements
    - wrt
        - Slides and their presentation
        - Web information
    - After lecture or during consultation hours
    - Or via email: naumann@hpi.uni-potsdam.de

# Textbook

A new first Edition textbook on

## Search Engines:
## Information Retrieval in Practice

Bruce Croft, *University of Massachusetts, Amherst*
Donald Metzler, *Yahoo Research*
Trevor Strohman, *Google*

Written by a leader in the field of information retrieval, *Search Engines: Information Retrieval in Practice* is designed to give undergraduate students the understanding and tools they need to evaluate, compare and modify search engines. The book covers the important issues in IR at a level appropriate for undergraduate computer science or computer engineering majors. Key mathematical models are included. The programming exercises in the book make extensive use of Galago, a Java-based open source search engine.

**Authors**

**Table of Contents**

**Sample Chapters**

A full draft of the textbook is available for review or for use in your information retrieval class.

For more information, please contact Matt Goldstein: matt.goldstein@aw.com

Estimated pub date is 2/15/09.

PEARSON
Addison Wesley

# Textbook



Search Engines
Information Retrieval
in Practice

W. BRUCE CROFT
DONALD METZLER
TREVOR STROHMAN

# Textbook

- 20 copies in library
- 75,99 € at amazon.de
  - Ouch, see
    http://www.newyorker.com/archive/2005/11/07/051107ta_talk_surowiecki
  - „ When professors decide which books to assign, the main consideration, they would say, is quality, not price, so any competition occurs on the basis of features rather than of cost. […] When price is no object, professors might as well choose the fanciest textbook around."
  - But: Free delivery…

# Other literature

- *Introduction to Information Retrieval*
  - Cambridge University Press, 2008.
  - Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze.
  - http://www-csli.stanford.edu/~hinrich/information-retrieval-book.html
- *Modern Information Retrieval*
  - Addison Wesley (27. Mai 1999)
  - Ricardo Baeza-Yates und Berthier Ribeiro-Neto

# Other literature - background

- The Google Story
    - http://www.amazon.de/Die-Google-Story-David-Vise/dp/3938017562/

# Other literature – Search Engine Optimization

16

- Which semester?

- HPI or IfI?

- Erasmus o.ä.?

- DB knowledge?

- Other relevant courses?
  - Semantic Web
  - Information Retrieval

- Your motivation?
  - Search engine optimization
  - Behind the scenes
  - Build your own search engine
  - Find a good job
  - Gain knowledge? Start research?

# Overview

- Introduction to team
- Organization
- Information Retrieval & Search Engines
- Overview of semester

# Search and Information Retrieval

- Search on the Web[1] is a daily activity for many people throughout the world

- Search and communication are most popular uses of the computer

- Applications involving search are everywhere

- The field of computer science that is most involved with R&D for search is *information retrieval (IR)*

[1] or is it <u>web</u>?

# Information Retrieval

*"Information retrieval is a field concerned with the structure, analysis, organization, storage, searching, and retrieval of information."* (Salton, 1968)

- General definition that can be applied to many types of information and search applications

  - □ Still appropriate after 40 years.

- Primary focus of IR since the 50s has been on *text* and *documents*

[1] or is it <u>web</u>?

http://www.cs.cornell.edu/Info/Department/Annual95/Faculty/Salton.html

# What is a Document?

- Examples:
  - Web pages, email, books, news stories, scholarly papers, text messages, Word™, Powerpoint™, PDF, forum postings, patents, IM sessions, etc.
- Common properties
  - Significant text content
  - Some structure ($\approx$ attributes in DB)
    - Papers: title, author, date
    - Email: subject, sender, destination, date

# Documents vs. Database Records

- Database records (or *tuples* in relational databases) are typically made up of well-defined fields (or *attributes*)
  - □ Bank records with account numbers, balances, names, addresses, social security numbers, dates of birth, etc.
- Easy to compare fields with well-defined semantics and data types to queries in order to find matches
  - □ Joins, selection predicates
- Text is more difficult, because unstructured

# Documents vs. Database Records

- Example bank database query
  - *Find records with balance > €50,000 in branches located in 14482 Potsdam.*
  - Matches easily found by comparison with field values of records
- Example search engine query
  - *bank scandals in western Germany*
  - This text must be compared to the text of *many, entire* news stories
    - ◇ Only "fields" might be *title* and *location*
- Defining the meaning of "balance" is much easier than defining "bank scandal".

# Comparing Text

- Comparing the query text to the document text and determining what is a good match is the <u>core issue</u> of information retrieval
- Exact matching of words is not enough
  - □ Many different ways to write the same thing in a "natural language" like English
    - ◇ Does a news story containing the text *"bank director in Potsdam steals funds"* match the query?
  - □ Some stories will be better matches than others
- Defining the meaning of a word, a sentence, a paragraph, or a story is more difficult than defining the meaning of a database field.

# Dimensions of IR

http://www.flickr.com/photos/garibaldi/3122956960/

- IR is more than just text, and more than just web search
  - □ although these are central
- People doing IR work with different media, different types of search applications, and different tasks

- Three dimensions of IR
  - □ Content
  - □ Applications
  - □ Tasks
- New applications increasingly involve new media
  - □ Video, photos, music, speech
  - □ Scanned documents (for legal purposes)
- Like text, content is difficult to describe and compare
  - □ Text may be used to represent them (e.g. tags)
- IR approaches to search and evaluation are appropriate

Tags

- germany
- 2008
- sanssouci
- brandenburg
- potsdam
- architectute
- castle
- garden
- clouds
- sky
- hdr
- 1xp
- photomatix
- lightroom
- gimp
- garibaldi
- column
- yellow
- autumn
- klausberg

# The Content Dimension

- New applications increasingly involve new media
    - □ Video, photos, music, speech
    - □ Scanned documents (for legal purposes)
- Like text, content is difficult to describe and compare
    - □ Text may be used to represent them (e.g. tags)
- IR approaches to search and evaluation are appropriate

http://www.flickr.com/photos/garibaldi/3122956960/



Tags

- germany
- 2008
- sanssouci
- brandenburg
- potsdam
- architectute
- castle
- garden
- clouds
- sky
- hdr
- 1xp
- photomatix
- lightroom
- gimp
- garibaldi
- column
- yellow
- autumn
- klausberg

# The Application Dimension

- Web search
    - Most common
- Vertical search
    - Restricted domain/topic
    - Books, movies, suppliers
- Enterprise search
    - Corporate intranet
    - Databases, emails, web pages, documentation, code, wikis, tags, directories, presentations, spreadsheets
- Desktop search
    - Personal enterprise search
    - See above plus recent web pages
- P2P search
    - No centralized control
    - File sharing, shared locality
- Literature search
- Forum search
- …

# The Task Dimension

- **User queries / ad-hoc search**
  - Range of query enormous, not prespecified
- **Filtering**
  - Given a profile (interests), notify about interesting news stories
  - Identify relevant user profiles for a new document
- **Classification / categorization**
  - Automatically assign text to one or more classes of a given set.
  - Identify relevant labels for documents
- **Question answering**
  - Similar to search
  - Automatically answer a question posed in natural language
  - Provide concrete answer, not list of documents.

**I ♥ INFERRET**

**Answers.com™**

How high is mt everest?    [ Ask ]

Recent questions:

What was the first civilization in America?

What was Houdini's most

Mt Everest is about twenty-nine thousand, five     hundred feet above sea level, making it the world's tallest mountain above sea level

👍 💬

http://amos.indiana.edu/library/scripts/mileshigh.html

**SHORT ANSWERS** ‹less / n

Answers 1-5

29035 FEET
8848
—
8850
AT 29035

- **Relevance**
  - □ A relevant document contains the information a user was looking for when he/she submitted the query.
- **Evaluation**
  - □ How well does the ranking meet the expectation of the user.
- **Users and information needs**
  - □ Users of a search engine are the ultimate judges of quality.

## Dead Search Engines   http://www.searchengineshowdown.com/reviews/

These search engines used to offer their own database or unique search features. They have all abandoned their position in search, although they still may have some kind of search functionality. The linked reviews reflect how these search engines used to work.

- AlltheWeb [Switched to Yahoo! database in March 2004]
- AltaVista [Switched to Yahoo! database in March 2004]
- Britannica Directory [some Web sites still included in the commercial Britannica, but not in the free version]
- Deja.com [Defunct Usenet search, bought by Google and became Google Groups]
- Direct Hit [Defunct, redirecting to Teoma]
- Excite [Defunct as a separate database. Now uses an InfoSpace meta search]
- Excite News (NewsTracker) [Defunct]
- Flipper [Hidden Web databases from Quigo, defunct by Fall 2003]
- Go [Defunct as a separate database, took over Infoseek, switched to Overture, then to Google]
- Go (Infoseek) News [Defunct]
- Infoseek [Defunct as a separate database, bought by Disney for Go, then abandoned in favor of Overture]
- HotBot [Dropped Inktomi database in early 2005, now only a multi-search of Google and Ask Jeeves]
- InvisibleWeb.com [a hidden Web directory, defunct by 2003]
- iWon [Old Inktomi version defunct. Now uses Google "sponsored" ads and Web and image databases]
- LookSmart [Directory
- Lycos [Switched to Yahoo!/Inktomi database in April 2004 and Ask Jeeves in 2005.]
- Magellan [Dead, redirects to WebCrawler]
- MessageKing [Defunct Web forum search engine as of Fall 2003]
- MSN Search [predecessor of Live Search]
- NBCi (formerly Snap) [Defunct, now uses metasearch engine Dogpile]
- NBCi Live Directory (formerly Snap) [Defunct directory]
- Northern Light [Defunct as a Web search engine as of 2002.]
- Northern Light Current News [Dead. Updates ceased as of Feb. 28, 2003.]
- Openfind [Under "reconstruction" as of 2003]
- Teoma [Dead, technology bought and now used by Ask.com]
- WebCrawler [Defunct as a separate database. Now uses an InfoSpace meta search]
- WebTop [Dead]
- WiseNut [Died in 2007]

# Relevance

- What is it?
- Simple (and simplistic) definition: A relevant document contains the information that a person was looking for when they submitted a query to the search engine
- Many factors influence a person's decision about what is relevant
  - Task at hand, context, novelty, style
- *Topical relevance* (same topic)
  - "Storm in Potsdam last Sunday" is topically relevant to query "Wetterereignisse"…
- vs. *user relevance* (everything else)
  - … but might not be relevant to user because
    - ◇ Read it before
    - ◇ Is five years old
    - ◇ Is in a foreign language, etc.

- *Retrieval models* define a view of relevance
  - Formal representation of the process of matching a query and a document.
  - Simple text matching as in DBMS or UNIX `grep` is not sufficient: Vocabulary mismatch problem
- *Ranking algorithms* used in search engines are based on retrieval models
  - Produce ranked list of documents
  - Real-world search engines consider topical and user relevance
- Most models describe statistical properties of text rather than linguistic
  - i.e. counting simple text features such as words instead of parsing and analyzing the sentences
  - Statistical approach to text processing started with Hans Peter Luhn in the 50s
    - ◇ Statistical view of text only recently popular in Natural Language Processing (NLP)
  - Linguistic features can be part of a statistical model

http://www.libsci.sc.edu/bob/chemnet/chist10.htm

http://www.lunometer.com/

# Evaluation

- Experimental procedures and measures for comparing system output with user expectations
    - □ Originated in Cranfield experiments in the 60s
- IR evaluation methods now used in many fields
- Typically use *test collection* (corpus) of documents, queries, and relevance judgments
    - □ Most commonly used are TREC collections
- *Recall* and *precision* are two examples of <u>effectiveness</u> measures
    - □ Precision: Proportion of retrieved documents that are relevant
    - □ Recall: Proportion of relevant documents that are retrieved
        - ◇ Assumption: All relevant documents are known. Ouch!
- Weblog data and clickthrough data to evaluate retrieval models and search engines.

# Users and Information Needs

- Search evaluation is user-centered

- Keyword queries are often poor descriptions of actual information needs
  - Query for "cats" could mean places to buy cats or the musical.
  - Search queries (in particular one-word queries) are under-specified.

- Interaction and context are important for understanding user intent

- Query refinement techniques such as
  - *query expansion*
  - *query suggestion*
  - *relevance feedback*

- improve ranking

Query length



http://www.submitexpress.com/news/shownews.php?article=1183

# IR and Search Engines

- A **search engine** is the practical application of information retrieval techniques to large scale text collections

- Web search engines are best-known examples, but many others exist

  - Web search: Crawl terabyte of web pages, provide subsecond response times, millions of queries

  - Enterprise search: variety of sources, search, data mining / clustering

  - Desktop search: rapidly incorporate new documents, many types of documents, intuitive interface

  - MEDLINE, online medical literature search since 70s

  - *Open source* search engines are important for research and development

    - ◇ Lucene, Lemur/Indri, Galago

- Big issues include main IR issues but also some others…

# IR and Search Engines

## Information Retrieval

- Relevance: *Effective ranking*
- Evaluation: *Testing and measuring*
- Information needs: *User interaction*

**Additional**

## Search Engines

- Performance: *Efficient search and indexing*
- Incorporating new data: *Coverage and freshness*
- Scalability: *Growing with data and users*
- Adaptability: *Tuning for applications*
- Specific problems: *e.g. Spam*

# Performance

- Measuring and improving the efficiency of search
  - □ reducing *response time*
  - □ increasing *query throughput*
  - □ increasing *indexing speed*
- *Indexes* are data structures designed to improve search efficiency
  - □ Designing and implementing them are major issues for search engines

- The "collection" for most real applications is constantly changing in terms of updates, additions, deletions
  - □ e.g., Web pages
- Acquiring or "crawling" the documents is a major task
  - □ Typical measures are *coverage* (how much has been indexed)
  - □ and *recency/freshness* (how recently was it indexed).
- Updating the indexes while processing queries is also a design issue

# Scalability

- Making everything work with millions of users every day, and many terabytes of documents

- Distributed processing is essential

- But: Large ≠ scalable

  - Scale gracefully

- Google in 2006

  - > 25 billion pages

  - 400M queries/day

- Google in 2008

  - 1 trillion pages (1,000,000,000,000)

    - http://googleblog.blogspot.com/2008/07/we-knew-web-was-big.html

# Adaptability

- Changing and tuning search engine components
    - ranking algorithm
    - indexing strategy
    - interface for different applications
- Adapt to different requirements for different applications / users

# Spam

- For Web search, spam in all its forms is one of <u>the</u> major issues
- Affects the efficiency of search engines and, more seriously, the <u>effectiveness</u> of the results
- Many types of spam
  - □ e.g. spamdexing or term spam, link spam, "optimization"
  - □ http://en.wikipedia.org/wiki/Spamdexing
- New subfield called *adversarial IR*, since spammers are "adversaries" with different goals

**Spamdexing** (also known as **search spam** or **search engine spam**)[1] involves a number of methods, such as repeating unrelated phrases, to manipulate the relevancy or prominence of resources indexed by a search engine, in a manner inconsistent with the purpose of the indexing system.[2][3] Some consider it to be a part of search engine optimization, though there are many search engine optimization methods that improve the quality and appearance of the content of web sites and serve content useful to many users.[4] Search engines use a variety of algorithms to determine relevancy ranking. Some of these include determining whether the search term appears in the META keywords tag

http://en.wikipedia.org/wiki/Spamdexing

# Overview

- Introduction to team
- Organization
- Information Retrieval & Search Engines
- → Overview of semester

# Chapter 2
# Architecture of a Search Engine

- **Basic Building Blocks**
- **Indexing**
  - □ Text Acquisition
  - □ Text Transformation
  - □ Index Creation
- **Querying**
  - □ User Interaction
  - □ Ranking
  - □ Evaluation

# Chapter 3
# Crawls and Feeds

- Deciding what to search

- Crawling the Web

- Directory Crawling

- Document Feeds

- The Conversion Problem

- Storing the Documents

- Detecting Duplicates

- Removing Noise

# Chapter 4
# Processing Text

- From Words to Terms

- Text Statistics

- Document Parsing

- Document Structure and Markup

- Link Analysis

- Information Extraction

- Internationalization

| | |
|---|---|
| Total documents | 84,678 |
| Total word occurrences | 39,749,179 |
| Vocabulary size | 198,763 |
| Words occurring > 1000 times | 4,169 |
| Words occurring once | 70,064 |

# Chapter 5
# Ranking with Indexes

- Abstract Model of Ranking

- Inverted indexes

- Compression

- Auxiliary Structures

- Index Construction – Map/Reduce

- Query Processing

Fred's **Tropical Fish** Shop is the best place to find **tropical fish** at low, low prices. Whether you're looking for a little **fish** or a big **fish**, we've got what you need. We even have fake **seaweed** for your fishtank (and little **surfboards** too).

9.7 fish
4.2 tropical
22.1 tropical fish
8.2 seaweed
4.2 surfboards

**Topical Features**

14 incoming links
3 days since last update

tropical fish
**Query**

**Ranking Function**

24.5
**Document Score**

**Document**

**Quality Features**

# Chapter 6
# Queries and Interfaces

- Information Needs and Queries

- Query Transformation and Refinement

- Showing the Results

- Cross-Language Search

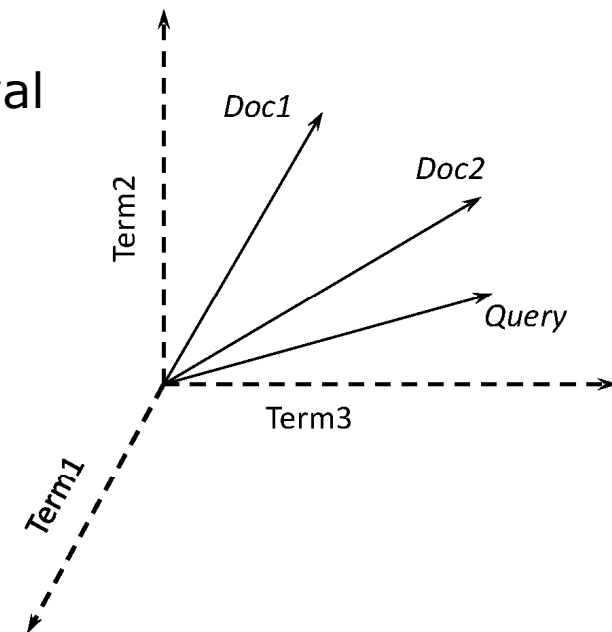| | | | | |
|---|---|---|---|---|
| 488941 britney spears | 29 britent spears | 9 brinttany spears | 5 brney spears | 3 britiy spears |
| 40134 brittany spears | 29 brittnany spears | 9 britanay spears | 5 broitney spears | 3 britmeny spears |
| 36315 brittney spears | 29 britttany spears | 9 britinany spears | 5 brotny spears | 3 britneeey spears |
| 24342 britany spears | 29 btiney spears | 9 britn spears | 5 bruteny spears | 3 britnehy spears |
| 7331 britny spears | 26 birttney spears | 9 britnew spears | 5 btiyney spears | 3 britnely spears |
| 6633 briteny spears | 26 breitney spears | 9 britneyn spears | 5 btrittney spears | 3 britnesy spears |
| 2696 britteny spears | 26 brinity spears | 9 britrney spears | 5 gritney spears | 3 britnetty spears |
| 1807 briney spears | 26 britenay spears | 9 brtiny spears | 5 spritney spears | 3 britnex spears |
| 1635 brittny spears | 26 britneyt spears | 9 brtittney spears | 4 bittny spears | 3 britneyxxx spears |
| 1479 brintey spears | 26 brittan spears | 9 brtny spears | 4 bnritney spears | 3 britnity spears |
| 1479 britanny spears | 26 brittne spears | 9 brytny spears | 4 brandy spears | 3 britntey spears |
| 1338 britiny spears | 26 btittany spears | 9 rbitney spears | 4 brbritney spears | 3 britnyey spears |
| 1211 britnet spears | 24 beitney spears | 8 birtiny spears | 4 breatiny spears | 3 britterny spears |
| 1096 britiney spears | 24 birteny spears | 8 bithney spears | 4 breetney spears | 3 brittneey spears |
| 991 britaney spears | 24 brightney spears | 8 brattany spears | 4 bretiney spears | 3 brittnney spears |
| 991 britnay spears | 24 brintiny spears | 8 breitny spears | 4 brfitney spears | 3 brittnyey spears |
| 811 brithney spears | 24 britanty spears | 8 breteny spears | 4 briattany spears | 3 brityen spears |
| 811 brtiney spears | 24 britenny spears | 8 brightny spears | 4 brieteny spears | 3 briytney spears |
| 664 birtney spears | 24 britini spears | 8 brintay spears | 4 briety spears | 3 brltney spears |
| 664 brintney spears | 24 britnwy spears | 8 brinttey spears | 4 briitny spears | 3 broteny spears |
| 664 briteney spears | 24 brittni spears | 8 briotney spears | 4 brittany spears | 3 brtaney spears |
| 601 bitney spears | 24 brittnie spears | 8 britanys spears | 4 brinie spears | 3 brtiiany spears |
| 601 brinty spears | 21 biritney spears | 8 britley spears | 4 brinteney spears | 3 brtinay spears |
| 544 brittaney spears | 21 birtany spears | 8 britneyb spears | 4 brintne spears | 3 brtinney spears |
| 544 brittnay spears | 21 biteny spears | 8 britnrey spears | 4 britaby spears | 3 brtitany spears |
| 364 britey spears | 21 bratney spears | 8 britnty spears | 4 britaey spears | 3 brtiteny spears |
| 364 brittiny spears | 21 britani spears | 8 brittner spears | 4 britainey spears | 3 brtnet spears |
| 329 brtney spears | 21 britanie spears | 8 brottany spears | 4 britinie spears | 3 brytiny spears |
| 269 bretney spears | 21 briteany spears | 7 baritney spears | 4 britinney spears | 3 btney spears |
| 269 britneys spears | 21 brittay spears | 7 birntey spears | 4 britmney spears | 3 drittney spears |
| 244 britne spears | 21 brittinay spears | 7 biteney spears | 4 britnear spears | 3 pretney spears |
| 244 brytney spears | 21 brtany spears | 7 bitiny spears | 4 britnel spears | 3 rbritney spears |
| 220 breatney spears | 21 brtiany spears | 7 breateny spears | 4 britneuy spears | 2 barittany spears |
| 220 britiany spears | 19 birney spears | 7 brianty spears | 4 britnewy spears | 2 bbbritney spears |
| 199 britrney spears | 19 brirtney spears | 7 brintye spears | 4 britnmey spears | 2 bbitney spears |
| 163 britnry spears | 19 britnaey spears | 7 britianny spears | 4 brittaby spears | 2 bbritny spears |
| 147 breatny spears | 19 britnee spears | 7 britly spears | 4 brittery spears | 2 bbrittany spears |
| 147 brittiney spears | 19 britony spears | 7 britnej spears | 4 britthey spears | 2 beitany spears |

# Chapter 7
# Retrieval Models

- Boolean Retrieval

- Vector Space Model

- Probabilistic Models

- Ranking based on Language Models

- Complex Queries and Combining Evidence

- Web Search

- Machine Learning and Information Retrieval
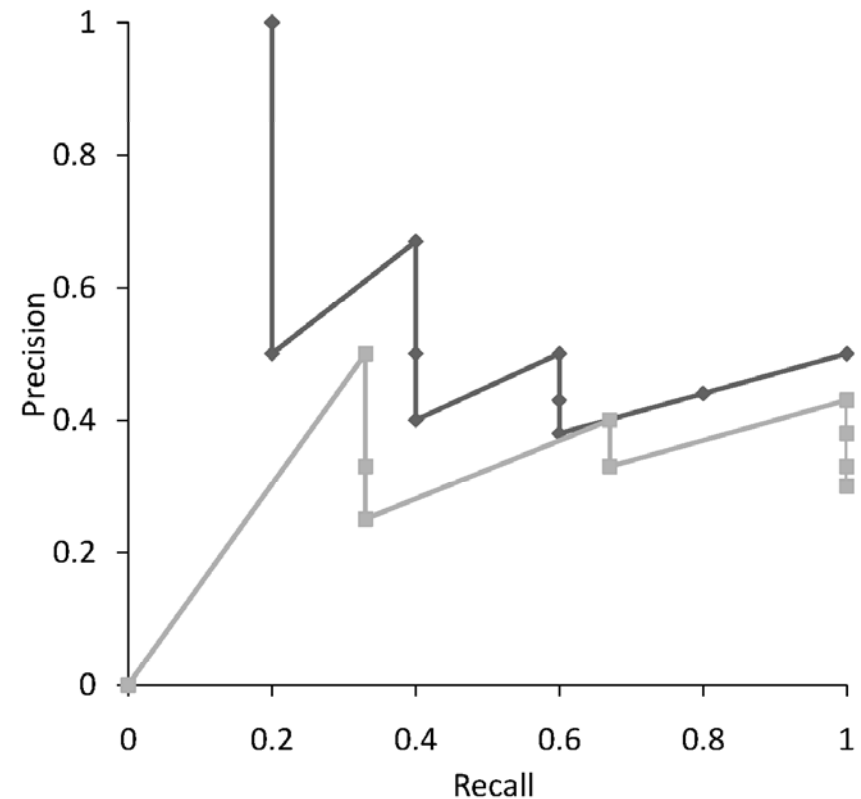
- Application-Based Models

# Chapter 8
# Evaluating Search Engines

- Motivation

- The Evaluation Corpus

- Logging

- Effectiveness Metrics

- Efficiency Metrics
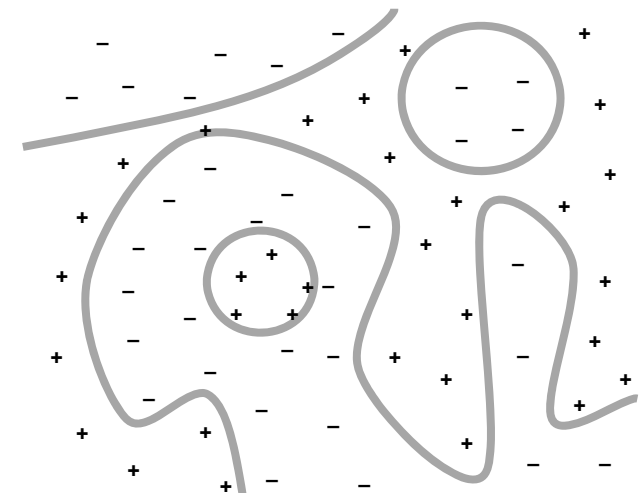
- Training, Testing, and Statistics

- The Bottom Line

# Chapter 9
# Classification and Clustering

- **Classification and Categorization**

  - Naïve Bayes

  - Support Vector Machines

  - Evaluation

  - Classifier and Feature Selection

  - Spam, Sentiment, and Online Advertising

- **Clustering**

  - Hierarchical and *K-Means Clustering*

  - *K Nearest Neighbor Clustering*

  - Evaluation

  - How to Choose K

  - Clustering and Search

# Chapter 10
# Social Search

- User Tags and Manual Indexing

- Searching With Communities

- Filtering and Recommending

- Document Filtering

- Personalization

- Peer-to-Peer and Metasearch

animals architecture art australia autumn baby band barcelona beach berlin birthday black blackandwhite blue california cameraphone canada canon car cat chicago china christmas church city clouds color concert day dog england europe family festival film florida flower flowers food france friends fun garden germany girl graffiti green halloween hawaii holiday home house india ireland italy japan july kids lake landscape light live london macro me mexico music nature new newyork night nikon nyc ocean paris park party people portrait red river rock sanfrancisco scotland sea seattle show sky snow spain spring street summer sunset taiwan texas thailand tokyo toronto travel tree trees trip uk usa vacation washington water wedding

# Chapter 11
# Beyond Bag of Words

- Feature-Based Retrieval Models

- Term Dependence Models

- Structure Revisited

- Longer Questions, Better Answers

- Words, Pictures, and Music

- One Search Fits All?



people, pool, swimmers, water

cars, formula, tracks, wall

clouds, jet, plane, sky

fox, forest, river, water

# Questions, wishes, …

- Now, or …
- Office: A.1-13
- Consultations: Tuesdays 15-16 Uhr
  or by arrangement
- Email: naumann@hpi.uni-potsdam.de
- Phone: (0331) 5509 280

## The end.