# Maschinelle Sprachverarbeitung

Introduction to Information Retrieval

Ulf Leser

# SHK Stelle bei WBI (~15h/Monat, 2 Jahre)

- Aufgabenbereich
  - Mitarbeit in einem Forschungsprojekt zur Analyse von Tabellen in wissenschaftlichen Texten
  - Erforschung und Entwicklung automatischer Verfahren zur Umstrukturierung tabellarischer Information
  - Entwicklung von Lernverfahren zur Ähnlichkeitsbewertung tabellarischer Daten

- Voraussetzungen
  - Studium der Informatik oder eines angrenzenden Bereichs
  - Sehr gute Kenntnisse in der Programmierung, vorzugsweise in Java oder Python
  - Gute Englischkenntnisse, mündlich und schriftlich
  - Kenntnisse in Verfahren der statistischen Sprachverarbeitung, der statistischen Datenanalyse, und des Maschinellen Lernens

# Content of this Lecture

- **What is Information Retrieval**
- Documents & Queries
- Text Preprocessing
- Evaluating IR Systems

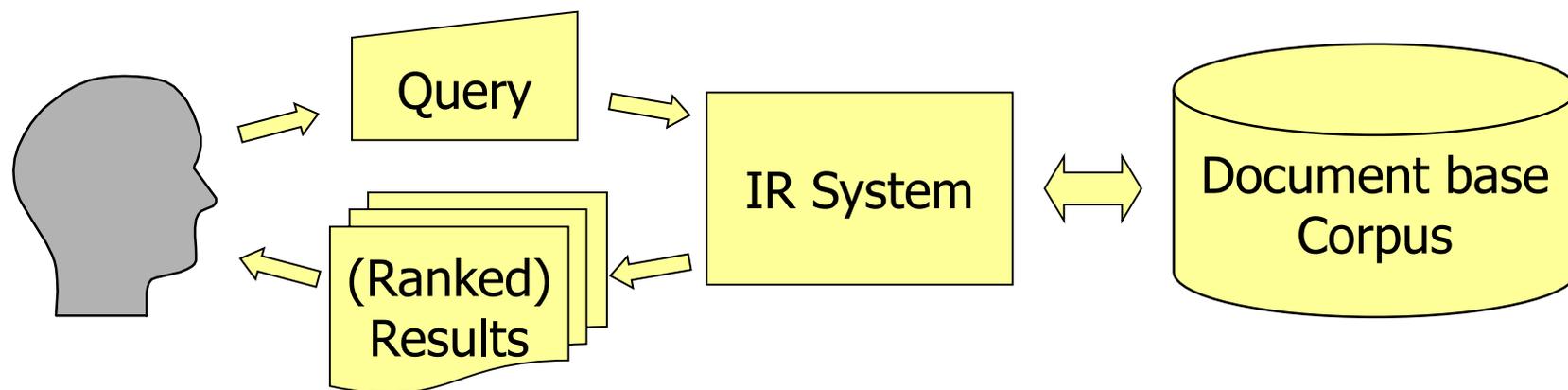# Information Retrieval (aka "Search")

- Naïve: Find all documents containing the following words
- Advanced: „Leading the user to those documents that will best enable him/her to satisfy his/her need for information"[Robertson 1981]
  - A user wants to know something
  - The user needs to tell the machine what he wants to know: query
  - Posing exact queries is difficult: room for interpretation
  - Machine interprets query to compute the (hopefully) best answer
  - Goodness of answer (relevance) depends on original intention of user, not on the query
  - "Leading": Sensible ranking of all potentially relevant docs

# The Problem

- Help user in quickly finding the requested information from a given set of documents
  - Documents: Corpus, library, collection, …
  - Quickly: Few queries, fast responses, simple interfaces, …
  - Requested: The "best-fitting" documents; the "right" passages; the most "relevant" content

Query

IR System
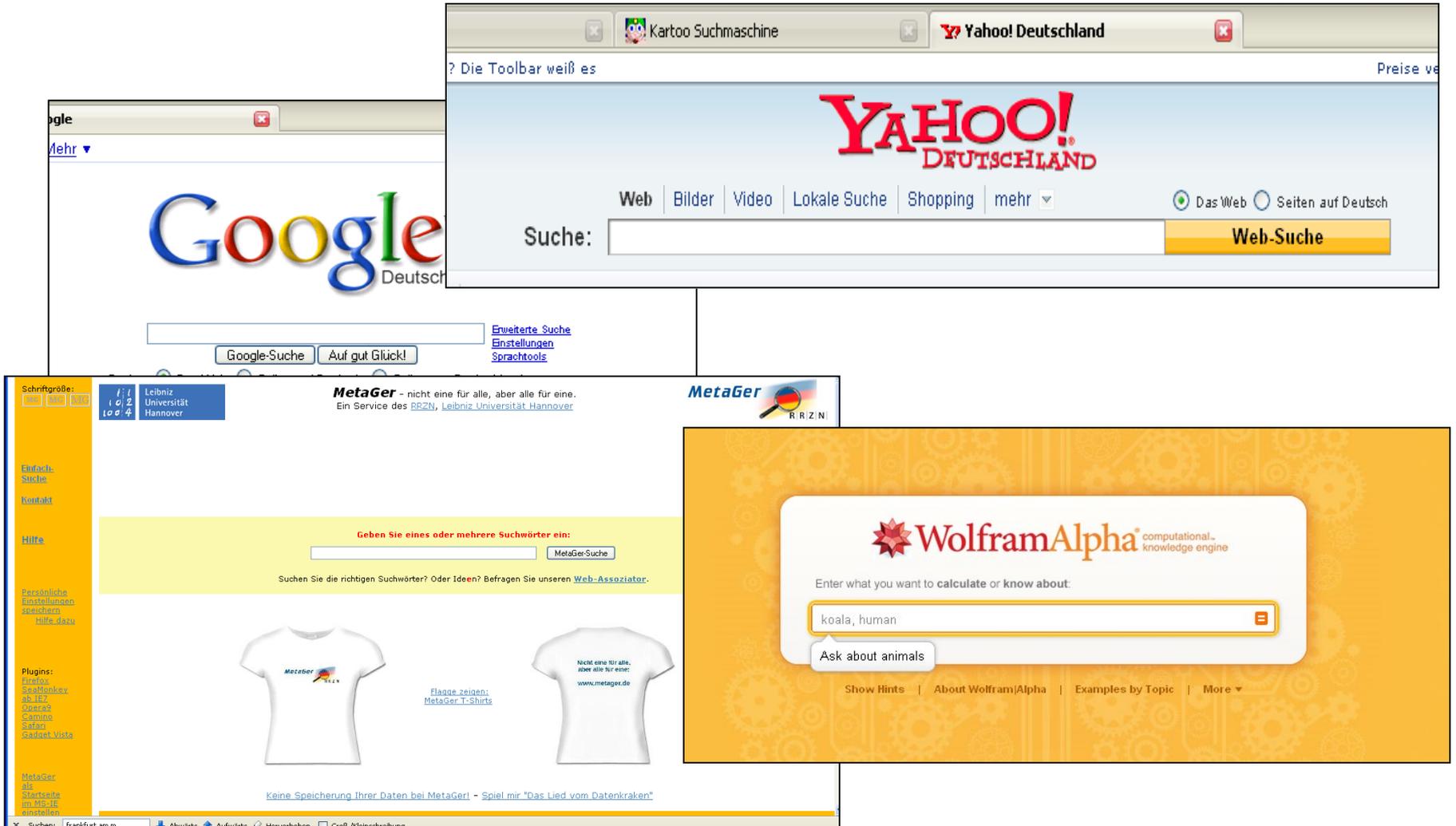
Document base Corpus

(Ranked) Results

# Why is it hard?

- Homonyms (context)
  - Fenster (Glas, Computer, Brief, …), Berlin (BRD, USA, …), Boden (Dach, Fussboden, Ende von etwas, …), …
- Synonyms
  - Computer, PC, Rechner, Server, Kiste, Bolide, Knoten, Desktop, …
- Specific queries (subfield Question Answering)
  - What was the score of Bayern München versus Stuttgart in the DFB Pokal finals in 1998? Who scored the first goal for Stuttgart?
  - How many hours of sunshine on average has a day in Crete in May?
- Typical web queries have 1,6 terms
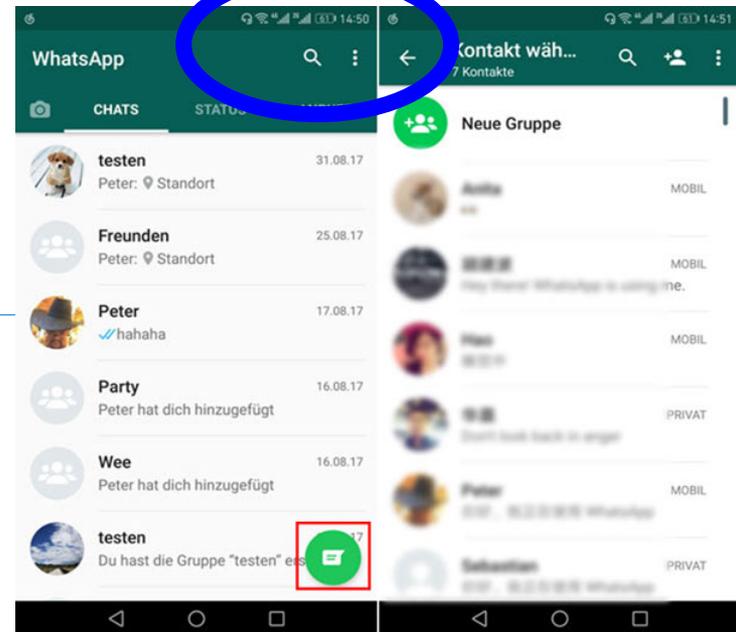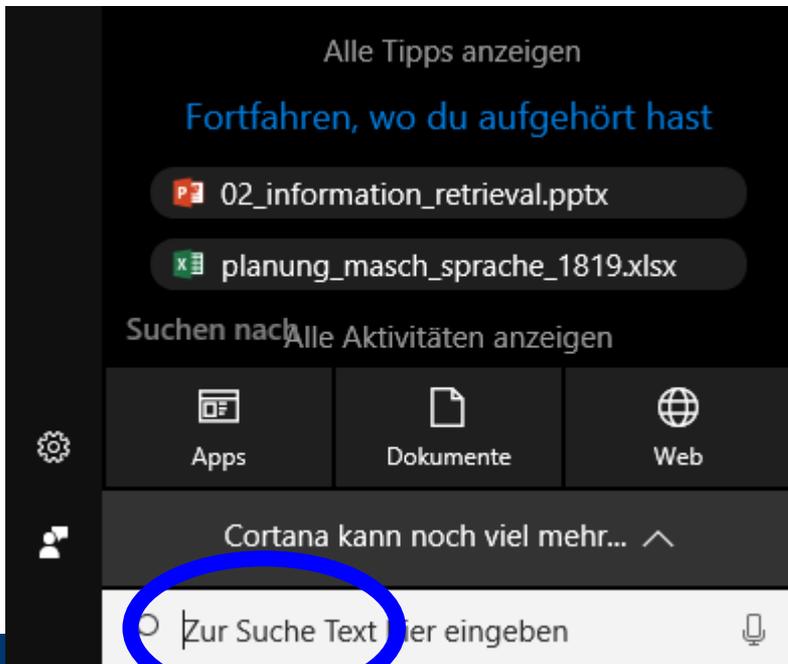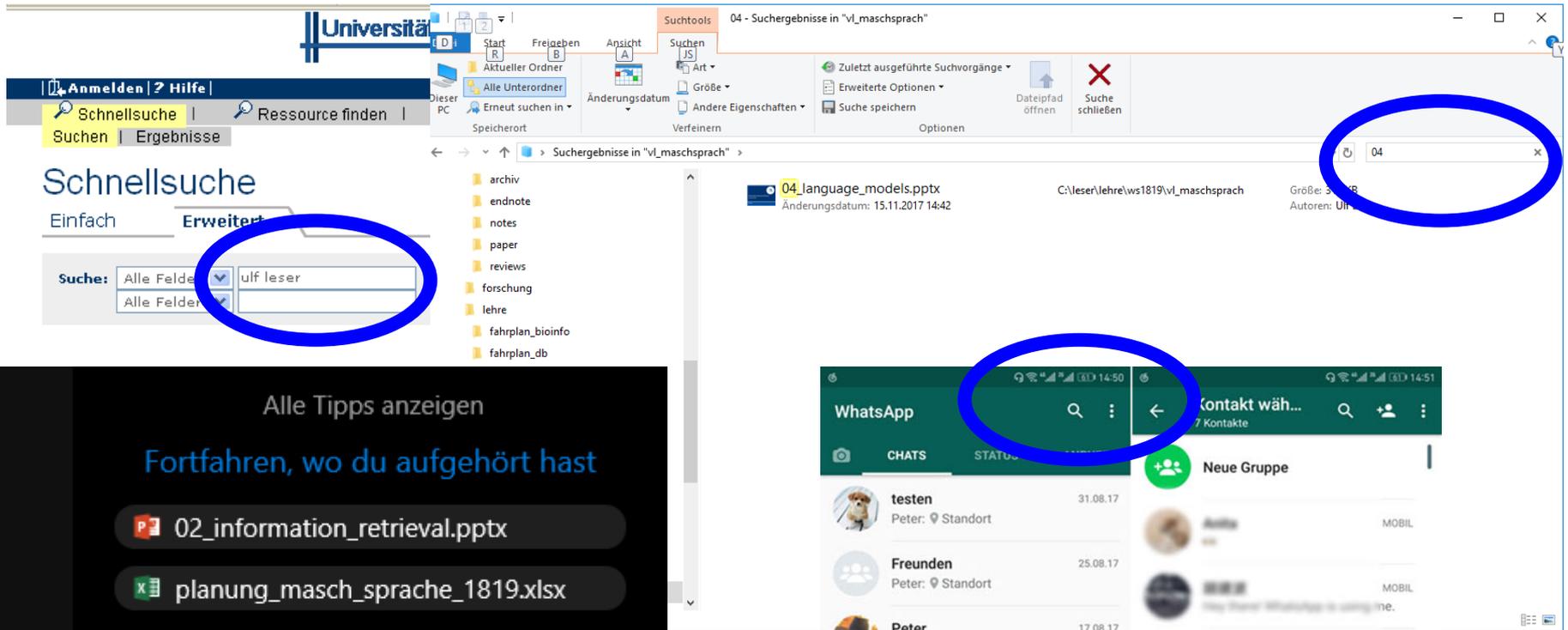- "Information broker" was (is?) a profession

# Quickly

- Time to execute a query
  - Indexing, parallelization, compression, …
- Time to answer the request (may involve multiple queries)
  - Understand request, find best matches
  - Success of search engines: Better results (and fast!)
  - Process-orientation: Exploit user feedback, query history, …
- Information overload
  - If the corpus is large, ranking is a must
  - Result summarization, result clustering
  - Different search modes: What's new? What's certain?

# Prominent Systems: Web Search Engines

# Other

# Content of this Lecture

- What is Information Retrieval
- Documents & Queries
- Text Preprocessing
- Evaluating IR Systems

# Document or Passage



Searching only metadata     Searching tokens within documents     Interpreting natural text

# Documents

- This lecture: Natural language text
- Might be grammatically correct (books, newspapers) or not (Blogs, Twitter, spoken language)
- May have structure (title, abstract, chapters, …) or not
- May have associated (explicit or in-text) metadata or not
- May be in many different languages or even mixed
  - Foreign characters
- May refer to other documents (hyperlinks)
- May have various formats (ASCII, PDF, DOC, XML, …)

# IR Queries

- Elements of IR-style queries
  - Keywords, phrases
  - Logical operations (AND, OR, NOT, …)
    - Web search: "-ulf +leser"
  - Structured queries on metadata (author=… AND title~ …)
- Documents as queries: Find documents similar to this one
- Query refinement based on previous results
  - Find documents matching the new query within the result set of the previous search
  - Use relevant answers from previous queries to create next query

# Searching with Metadata (PubMed/Medline)

# Content of this Lecture

- What is Information Retrieval
- Documents & Queries
- Text Preprocessing
  - Special characters and case
  - Sentence splitting
  - Tokenization
  - Stop words
  - Zipf's law
- Evaluating IR Systems

# Processing Pipeline



Frakes et al. 1992

# Logical View

- Definition
  - *The logical view of a document denotes its representation inside the IR system*

- Determines what users can address in a query
  - Only metadata, only title, only abstract, full text, phrases, stop words, special characters, …

- Creating the logical view involves transformations
  - Stemming, stop word removal
  - Transformation of special characters (Umlaute, Greek letters, …)
  - Removal of formatting information (HTML), tags (XML), …
  - Bag of words (BoW)
    - Arbitrary yet fixed order (e.g. sorted alphabetically)
    - See next lecture

# Definitions

- Definition
  - *A document as a sequence of sentences*
  - *A sentence is a sequence of tokens*
  - *A token is the smallest unit of text (words, numbers, …)*
  - *A concept is the mental representation of a "thing"*
  - *A term is a token or a set of tokens representing a concept*
    - *"San" is a token, but not a term*
    - *"San Francisco" has two tokens but is only one term*
    - *Dictionaries usually contain terms, not tokens*
  - *A homonym is a term representing multiple concepts*
  - *A synonym is a term representing a concept which may also be represented by other terms*
  - *A syn-set is a set of synonyms representing the same concept*
- "Word" can denote either a token or a term
- We will mostly make no difference between token and terms (sadly …)

# Format Conversion

> **ABSTRACT**
>
> New generation of e-commerce applications require data schemas that are constantly evolving and sparsely populated. The conventional horizontal row representation fails to meet these require-
>
> **1.1 Issues**
>
> In relational database systems, data objects are conventionally stored using a horizontal scheme. A data object is represented as a row of a table. There are as many columns in the table as the number of attributes the objects have. In trying to store all our

> New generation of e-commerce applications require data schemas In relational database systems, data objects are conventionally that are constantly evolving and sparsely populated. The conven- stored using a horizontal scheme. A data object is represented as tional horizontal row representation fails to meet these require- …

- Transform PDF, XML, DOC, …  into ASCII / UniCode
- Problems: Formatting instruction, special characters, formulas, tables, section headers, footnotes, …
- Web: Find the net content (no ads, navigation bars, …)
- Diplomacy: To what extend can one reconstruct the original document from its logical view?

# Case – A Difficult Case

- Should all text be converted to lower case letters?
- Advantages
  - Makes queries simpler
  - Decreases index size
  - Allows for some "fuzziness" in search
- Disadvantages
  - No abbreviations
  - Loss of important hints for sentence splitting
  - Loss of important hints for tokenization, NER, …
  - Loss of semantic info (proper names, Essen versus essen,…)
- Different impact in different languages (German / English)
- Often: Convert only after all other preprocessing steps

# Sentence Splitting

- Most linguistic analysis works on sentence level
- Sentences group together entities and statements
- Naive approach: Reg-Exp search for "[.?!;] "
  - (note the blank!)
  - Abbreviations
    - "C. Elegans is a worm which …"; "This does not hold for the U.S."
  - Errors (due to previous normalization steps)
    - "is not clear.Conclusions.We reported on …"
  - Proper names
    - ".NET is a technique for …"
  - Direct speech
    - "By saying "It is the economy, stupid!", B. Clinton meant that …"
  - …

# Tokenization

- Fundamental elements of all IR query languages are tokens
- Simple approach: search for „ " (blanks)
  - "A state-of-the-art Z-9 Firebird was purchased on 3/12/1995."
  - „SQL commands comprise SELECT … FROM … WHERE clauses; the latter may contain functions such as leftstr(STRING, INT)."
  - "This LCD-TV-Screen cost 3,100.99 USD."
  - "[Bis[1,2-cyclohexanedionedioximato(1-)-O]-[1,2-cyclohexanedione dioximato(2-)-O]methyl-borato(2-)-N,N0,N00,N000,N0000,N00000)-chlorotechnetium) belongs to a family of …"
- Typical approach (but many (domain-specific) variations)
  - Treat hyphens / parentheses as blanks
  - Remove "." (after sentence splitting)

# Stop Words

- Words that are so frequent that their removal (hopefully) does not change the meaning of a doc
  - English: Top-2: 10% of all tokens; Top6: 20%; Top-50: 50%
  - English (top-10; LOB corpus): the, of, and, to, a, in, that, is, was, it
  - German(top-100): aber, als, am, an, auch, auf, aus, bei, bin, …
- Consequences
  - Removing top-100 stop words reduces a positional index by ~40%
  - Hopefully increases precision due to less spurious hits
  - Makes many phrase queries impossible
- Variations
  - Remove top 10, 100, 1000, … words
  - Language-specific, domain-specific, corpus-specific

# Example

The children of obese and overweight parents have an increased  risk of obesity. Subjects with two obese parents are fatter in childhood and also show a stronger pattern of tracking from childhood to adulthood. As the prevalence of parental obesity increases in the general population the extent of child to adult tracking of BMI is likely to strengthen.

↓ 100 stop words

children obese overweight parents increased  risk obesity. Subjects obese parents fatter childhood show stronger pattern tracking childhood adulthood. prevalence parental obesity increases general population extent child adult tracking BMI likely strengthen.

↓ 10 000 stop words

obese overweight obesity obese fatter adulthood prevalence parental obesity BMI

# Zipf's Law

- Let f be the frequency of a word and r its rank in the list of all words sorted by frequency

- Zipf's law: $f \sim k/r$ for some constant k

- Example
  - Word ranks in Moby Dick
  - Good fit to Zipf's law
  - Some domain-dependency (whale)

- Fairly good approximation for most corpora

Source: http://searchengineland.com/the-long-tail-of-search-12198

# Content of this Lecture

- What is Information Retrieval
- Documents & Queries
- Text Preprocessing
- Evaluating IR Systems

# Evaluation: Binary Model

- Assume that for a given query q and many docs $d \in D$, somebody determines whether d is relevant for q or not
  - An expert? An average user?
- The IR systems returns all docs it thinks that are relevant
  - No ranking for now
- Let T be the set of all truly relevant docs, X the set of all returned docs: |T|=TP+FN, |X|=TP+FP

|  | Truth: relevant | Truth: not relevant |
|---|---|---|
| **IR: relevant** | True positives | False positives |
| **IR: not relevant** | False negatives | True negatives |

# Precision and Recall

- Precision = TP/(TP+FP)
  - What is the fraction of relevant answers in X?
- Recall = TP/(TP+FN)
  - What is the fraction of found answers in T?
- The perfect world

|  | Real: Positive | Real: Negative |
|---|---|---|
| IR: Positive | A | 0 |
| IR: Negative | 0 | B |

# Example

- Let |D| = 10.000, |X|=15, |T|=20, |X∩T|=9

| | Real: Positive | Real: Negative |
|---|---|---|
| IR: Positive | TP = 9 | FP = 6 |
| IR: Negative | FN = 11 | TN= 9.974 |

  – Precision = TP/(TP+FP) = 9/15 = 60%
  – Recall    = TP/(TP+FN) = 9/20 = 45%

- Assume another result: |X|=10, |X∩T|=7

| | Real: Positive | Real: Negative |
|---|---|---|
| IR: Positive | TP = 7 | FP = 3 |
| IR: Negative | FN = 13 | |

  – Precision: 70%, recall = 35%

# A Different View

Retrieved

Relevant

Excellent precision,
terrible recall

Terrible precision,
terrible recall

Terrible precision,
excellent recall

Excellent precision,
Excellent recall

# Trade-off
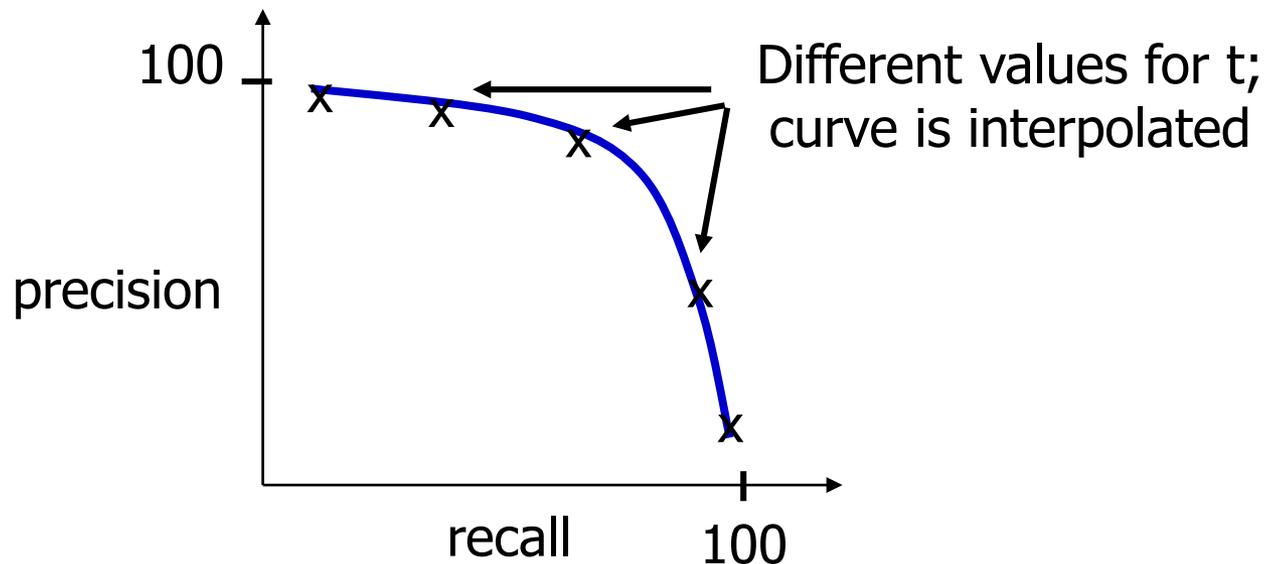
- Inherent Trade-off between precision and recall
- Example
  - Think VSM with a threshold t to enforce a binary decision
  - Assume that docs with high relevance score are most likely also relevant for the user
  - Increase t:       Less results, most of them very likely relevant
                      Precision increases, recall drops
                      Set t=1: P~100%, R=?
  - Decrease t:       More results, some might be wrong
                      Precision drops, recall increases
                      Set t=0: P=?, R=100%

# Precision / Recall Curve

- Sliding the threshold t gives a curve
  - Similar to Receiver-Operating-Characteristic (ROC)



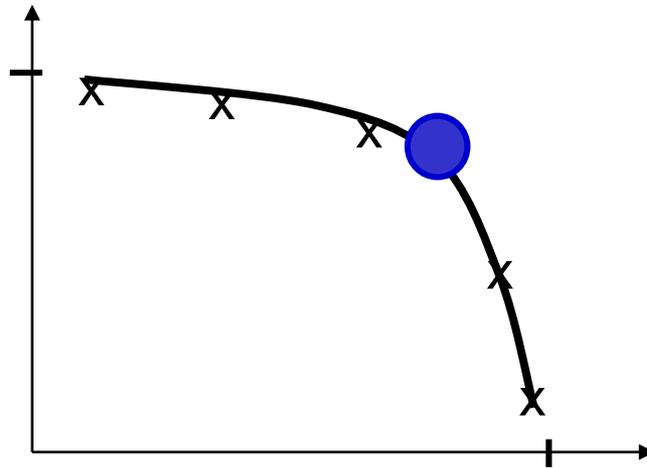Different values for t; curve is interpolated

100

precision

recall          100

# Accuracy

|  | Truth: relevant | Truth: not relevant |
|---|---|---|
| **IR: relevant** | TP | FP |
| **IR: not relevant** | FN | TN |

- Sometimes, one wants one measure instead of two
  - E.g. to rank different IR-systems
- Accuracy= (TP+TN) / (TP+FP+FN+TN)
  - Which percentage of the systems decision were correct?
  - Makes only sense with small corpora and large result set
  - Typically in IR, TN >>> TP+FP+FN
  - Thus, accuracy is always good (~0,9999995)
- Used in problems with balanced sets of TN / TP

# F-Measure

- F-Measure = 2*P*R / (P+R)
  - F-Measure is harmonic mean between precision and recall
  - Favors balanced P/R values



- Alternative: Area-under-the-curve (AUC)
  - AUC doesn't help in finding the best t

# From user/query to users/queries

- What if we have many queries?
  - Evaluation should always use a range of different queries
  - Compute average P/R values over all queries
  - Of course, mean and stddev are also important
- What if we have different users?
  - Different users may have different thoughts about what is relevant
  - Leads to different gold standards
  - Compute inter-annotator agreement as upper bound
- Who can judge millions of docs?
  - Evaluate on small gold standard corpus
    - But: Extrapolation is difficult: Are the properties of application/corpus really equal to properties of GS?
  - Use implicit feedback, e.g. click-through rates in top-K results

# Self Assessment

- Give a definition of „Information retrieval"
- How is information retrieval different from database query evaluation?
- What are possible types of answers to a IR query?
- List 5 important steps in document preprocessing and their expected impact on precision and recall
- Give a definition of recall, precision, and accuracy
- What is the difference between micro and macro average
- Which preprocessing steps would be affected if you work with a multi language corpus?