IR4IP Tutorial IPI Confex, March 2009 Erik Graf and Thomas Roelleke

N RETRIEN

Introduction IPI Confex, March 2009

Thomas Roelleke Queen Mary University of London

INFORMATION SERVICES		
1	1/18 Structure of IB4IP Tutorial	2/18
	Topics, Issues, and "Problems" in IR Foundations Historical Notes Web or Enterprise Search? Summary	
Outline	Structure	
<ol> <li>Structure of IR4IP Tutorial</li> <li>Topics, Issues, and "Problems" in IR</li> </ol>	1200 - 1330: 90 mins, "six" 15 mins slots.	
<ul> <li>Foundations</li> <li>Historical Natas</li> </ul>	<ul> <li>Indexing (EG): 15 mins</li> <li>Retrieval Models (TR): 15 mins</li> </ul>	
<ul><li>Instorical Notes</li><li>Web or Enterprise Search?</li></ul>	<ul> <li>Interaction (EG): 15 mins</li> <li>Evaluation (TR): 15 mins</li> </ul>	
6 Summary		

Structure of IR4IP Tutorial Topics, Issues, and "Problems" in IR Foundations **Historical Notes** Web or Enterprise Search? Summary

What does Wikipedia say? **IR Herbstschule 2008** Empty Answer and Many Answer "Problem"

Structure of IR4IP Tutorial Topics, Issues, and "Problems" in IR Foundations **Historical Notes** Web or Enterprise Search? Summary

10 Issues in IR

What does Wikipedia say? 10 Issues in IR **IR Herbstschule 2008** Empty Answer and Many Answer "Problem"

### What does Wikipedia say?

#### http://en.wikipedia.org/wiki/Information\_retrieval Retrieval models (ranking functions) Your continued donations keep Wikipedia running! 🚨 Log in / create account article discussion edit this page history Text processing ("Indexing"): NLP / understanding W 20 28 Help shape the future of Wikipedia. Please participate in our survey of readers and 5 N contributors! (More information) Interactivity すつ祖 22 Information retrieval Efficiency: Compression, parallel IR WIKIPEDIA From Wikipedia, the free encyclopedia Distributed IR (data fusion, meta retrieval) The Free Encyclopedia Information retrieval (IR) is the science of searching for documents, for information navigation Multimedia: image, video, sound within documents and for metadata about documents, as well as that of searching 6 Main page relational databases and the World Wide Web. There is overlap in the usage of the terms Contents data retrieval, document retrieval, information retrieval, and text retrieval, but each also has 7 Evaluation Featured content its own body of literature, theory, praxis and technologies. IR is interdisciplinary, based on Current events Web retrieval (link analysis) computer science, mathematics, library science, information science, information Random article architecture, cognitive psychology, linguistics, statistics and physics. search Cross-lingual IR Automated information retrieval systems are used to reduce what has been called "information overload". Many universities and public libraries use IR systems to provide Digital libraries: IR application Go Search 10 access to books, journals and other documents. Web search engines are the most visible IR applications. interaction see http://www.dlib.org/dlib/november95/11croft.html About Wikipedia Contents [hide] Community portal 5/18 Structure of IR4IP Tutorial Structure of IR4IP Tutorial Topics, Issues, and "Problems" in IR What does Wikipedia say? Topics, Issues, and "Problems" in IR What does Wikipedia say? Foundations Foundations 10 Issues in IR Historical Notes **IR Herbstschule 2008 Historical Notes IR Herbstschule 2008** Web or Enterprise Search? Empty Answer and Many Answer "Problem' Web or Enterprise Search? Empty Answer and Many Answer "Problem" Summary Summary Empty Answer and Many Answer "Problem"

7/18

### **IR Herbstschule 2008**

IR for PhD and post-doc researchers.

http://www.dagstuhl.de/en/program/calendar/evhp/?semnr=08402



- "Empty Answer Problem":
  - Free-text: Find a web page that offers boats that are fast AND comfortable AND NOT EXPENSIVE
  - SQL DB: SELECT \* FROM properties WHERE price < 200k AND bedrooms > 3 AND location LIKE 'London':

Query is too narrow, too specific.

- "Many Answer Problem":
  - SELECT \* FROM properties WHERE price < 200k OR</li> bedrooms > 3 OR location LIKE 'London';

Query is too general, too exhaustive.



- retrieval based on the distance of document and query vectors
- Probabilistic justification of what ranking is? P(r|d,q), the probability of relevance.

Illustration: Present a document-guery pair to several users, and each user assesses the relevance.

After ... steps: A term/keyword is GOOD, if it occurs more often in relevant documents than in non-relevant documents.

Query

Retrieval function

Retrieved documents

Document representation

Structure of IR4IP Tutorial Topics, Issues, and "Problems" in IR Foundations Historical Notes Web or Enterprise Search? Summary

The 60s/70s **The 80s** Mid 90s: The web Late 90s: Language modelling Structure of IR4IP Tutorial Topics, Issues, and "Problems" in IR Foundations Historical Notes Web or Enterprise Search? Summary

The 60s/70s The 80s Mid 90s: The web Late 90s: Language modelling

### The 80s

### Mid 90s: The web

- SQL databases become widely available ... text processing?
- VSM, probability of relevance, all interesting, but not good enough?
- Therefore, view IR as "logical implication" ... Rijsbergen's P(d 
  ightarrow q)

- The web. Keyword-based (content-based) retrieval alone not good enough.
- Pagerank ... Google ... A page is good if it is "popular".
- A popular page is an authority?
- Combine content-based (keyword-based) ranking with pagerank.
- The page rank (popularity, authority) is independent of the query!

	13/18		14/18
Structure of IR4IP Tutorial Topics, Issues, and "Problems" in IR Foundations Historical Notes Web or Enterprise Search? Summary	The 60s/70s The 80s Mid 90s: The web Late 90s: Language modelling	Structure of IR4IP Tutorial Topics, Issues, and "Problems" in IR Foundations Historical Notes Web or Enterprise Search? Summary	
Late 90s: Language modelling		Web/page search OR Ente	erprise search OR Semantic
		web?	

- VSM, probabilistic model, logical model, all interesting, but how about language modelling?
- Given a document and a collection, there is a probability that the query is generated from those *TWO* "models"!

- Web search: surface/horizontal search
- Enterprise search: "deep" search, semantic search, vertical search
- Semantic web, webDB: explore/exploit the web similar to what is known for enterprise search

Structure of IR4IP Tutorial Topics, Issues, and "Problems" in IR Foundations Historical Notes Web or Enterprise Search? Summary

## Summary

## Questions?

- Information retrieval: not just "document" retrieval
- An information need is translated into a query: "loss of information"
- Hypertext/Web: Content (words/terms) + Links (links, popularity, authority)

## Thank you!

17/18



Introduction Full term indices Advanced Index structures Conclusion

Concept Very Short History of Indexing Introduction Full term indices Advanced Index structures Conclusion

Concept Very Short History of Indexing

## **Overview of Indexing Process**

## **Overview of Indexing Process**



Index creation

Introduction Full term indices Advanced Index structures Conclusion

### Index creation



### Index creation

#### Term Processing

- Tokenization: Extract the terms from a document. Removal of tags, punctuation.
- Stop-wording: Remove terms with very high document frequencies (e.g. 'of', 'the'). In the English language these are responsible for approx. 30% of term volume.

Index creation

Stop word lists are getting shorter; or often no stop-wording is applied.

• Stemming: Collapsing morphological variants of words (i.e 'Patent', 'Patents', 'patenting' is stemmed to 'patent'). Eases guerying and reduces index size.

9/18 10/18Introduction Introduction Full term indices Full term indices Index creation Index creation Advanced Index structures Advanced Index structures Conclusion Conclusion Inverted Index

### **Direct Index**

_	
DOC1	 Term1, Term3, Term4, Term5,
DOC2	 Term1, Term3, Term4, Term5,
DOC3	 Term1, Term3, Term4, Term5,
DOC4	 Term1, Term3, Term4, Term5,
DOC5	 Term1, Term3, Term4, Term5,

- Optimized file system for retrieval. Mimics a Unix 'grep' or Windows98 search.
- Scales badly with respect to number of documents.



- Default index structure in Information Retrieval.
- Computationally very efficient. Scales well.

	Intro	auction
Full	term	indices
Advanced Ind	ex str	uctures
	Con	clusion

#### Positional Index Field Index

Introduction Full term indices Advanced Index structures Conclusion

Positional Index Field Index Virtual Fields

## **Positional Index**

## Field Index



- Stores the position of term in addition to the posting and frequency
- Allows for word-order- ,phrasal- , offset-querying.
- Increases the size of the index by 2-10 times (depending on the average document length)

		12	United States Patent	the Passer: Nac. 405 (1) Date of Palence	C514,142 B1		
	/	194 195 177 127 127 127 127 127 127 127 127 127	A CONTRACTOR OF A CONTRACTOR A C	1423 1 202 202, 2023 1 202 202, 2023 1 202 202, 2023 2 202 202, 2023 2 202 202, 2023 2 202 202, 2024 2 20, 2024 2 20, 2024 2 20	in a second seco	_	
	ītle	Ab	stract	De	▼ script		► Laims
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings
Term	Postings	Term	Postings	Term	Postings	Tern	n Postings

	13/18			14/18
Introduction Full term indices Advanced Index structures Conclusion		Introduction Full term indices Advanced Index structures Conclusion	Positional Index Field Index Virtual Fields	
Field Index		Virtual fields		

- Enables part specific searching
- Allows assigning weight to different parts or aspects of documents (Robertson BM25F)
- A form of partitioning documents.

#### Virtual fields

- Enrich document with document-external data
- Document priors (Web: Number of inlinks, URL length, URL text)
- Another example anchor text:
  - Link: 'http://www.acm.org/sigir/'
  - Text: 'ACM SIGIR: Information Retrieval Special Interest Group'

Introduction Full term indices Advanced Index structures Conclusion

## Summary

#### 'Trends' in document representations

- Representations resemble documents more closely
  - Full text
  - Positional
- Exploitation of semantic 'hints'
  - Structure
  - Style
- Enrichment of document with meta-data
  - Anchor text
  - Document priors (inlinks,URL, etc ...)

# Thank you!

Questions?

17/18

IR4IP Tutorial IPI Confex, March 2009 Erik Graf and Thomas Roelleke Retrieval Models IPI Confex, March 2009

Thomas Roelleke Queen Mary University of London



Introduction **Retrieval Models Related Models/Tasks Relationships between Retrieval Models** Summary

## Time frame

60s/70s: Vector-space model (VSM); TF-IDF (term	Term sailing boats sailing	Docld doc1 doc1 doc2	$tf(t, d)$ term frequency $df(t)$ document frequency $N_D$ number of documents
probabilistic (binary independence) retrieval model	boats sailing	doc2 doc2	tf(sailing, doc2) = 2
(BIR)	sailing	doc3	df(sailing) = 4
Early/mid 90s: Poisson, BM25 (best-match version 25)	east	doc3	$N_D = 5$
Mid 90s: Page-rank (Authority-based retrieval a-la Google) Late 90s: Language Modelling	coast sailing boats	doc3 doc4 doc5	$idf(sailing) = -\log \frac{4}{5} = \dots$

	5/18		6/1
Introduction Retrieval Models Related Models/Tasks Relationships between Retrieval Models Summary	Vector-Space Model (VSM) TF-IDF Binary Independence Retrieval (BIR) Model Poisson Model BM25 Formula Language Modelling (LM)	Introduction Retrieval Models Related Models/Tasks Relationships between Retrieval Models Summary	Vector-Space Model (VSM) TF-IDF Binary Independence Retrieval (BIR) Model Poisson Model BM25 Formula Language Modelling (LM)
Vector-Space Model (VSM	)	TF-IDF	





## TF: Term Frequency of TERM in DOCUMENT IDF: Inverse Document Frequency of TERM in COLLECTION

tf(t, d) := count term occurrence within DOCUMENT

idf(t) := count documents in COLLECTION

IDF of frequent term is small, IDF of rare term is large. Reflects the searcher trying to find terms that are rare/discriminative overall, but frequent in the requested document.

Retrieval Models Related Models/Tasks Relationships between Retrieval Models Summary Vector-Space Model (VSM) TF-IDF Binary Independence Retrieval (BIR) Model Poisson Model BM25 Formula Language Modelling (LM)

### Binary Independence Retrieval (BIR) Model

Introduction

Introduction Retrieval Models Related Models/Tasks Relationships between Retrieval Models Summary Vector-Space Model (VSM) TF-IDF Binary Independence Retrieval (BIR) Model Poisson Model BM25 Formula Language Modelling (LM)

### Poisson Model

- Established mid/end 70s ([Robertson and Sparck Jones, 1976, Croft and Harper, 1979])
- Terms are "good": if they are frequent in relevant documents and rare in non-relevant documents
- Terms are "poor": if they are frequent in non-relevant documents and rare in relevant documents
- Famous formula (BIR term weight):

$$\mathsf{bir}(t,r,\bar{r}) := \frac{P(t|r)}{P(t|\bar{r})}$$

- Probability that an event occurs k times given that in average it occurs λ times
- Example: Probability that 4 of 7 days are sunny, knowing that in average every second day was sunny in the past (one month: 15/30, 10 years: 1800/3600).



- established since early/mid 90s
- combines a special TF component with the BIR term weight
- considers document length
- is mathematically

$$\sum_{r \in d \cap q} \mathsf{bir}(t, r, \overline{r}) \cdot \frac{\mathsf{tf}_d}{\mathsf{tf}_d + k_1} \cdot \frac{\mathsf{tf}_q}{\mathsf{tf}_q + k_3} + k_2 \cdot \mathsf{qI} \cdot \frac{\mathsf{avgdI} - \mathsf{dI}}{\mathsf{avgdI} + \mathsf{dI}}$$

 $k_1$ : pivoted document length ...

- THE alternative to BM25/TF-IDF?
- Established late 90s
- Derived from P(q|d), i.e. the probability that document d and the background model (the collection) generate the query.
- Example: ...

Introduction Retrieval Models Related Models/Tasks Relationships between Retrieval Models Summary

Classification-oriented Models Web/Link-based Retrieval Models Introduction Retrieval Models Related Models/Tasks Relationships between Retrieval Models Summary

Classification-oriented Models Web/Link-based Retrieval Models

## **Classification-oriented Models**

## Web/Link-based Retrieval Models

- Bayesian classifier
- Support-vector machine (SVM)
- Duality to ad-hoc retrieval: Retrieve classes for an incoming (new) document

- Idea: A page that is referenced by many "good" pages is a "good" page. Note the recursive usage of "good" ([Brin and Page, 1998]).
- Authority principle.
- Apparently the break-through for Google late 90s.
- TF boosting: Propagate the anchor text terms to the referenced object; multimedia/image retrieval.



- The content-oriented models (TF-IDF, BM25, LM) are combined with link-based models (e.g. term propagation).
- All retrieval models try to optimise the ranking.
- Can one know in advance which model is best for which query?
- Is a combination of models useful?
- Can a system learn a model? Learn when to use which model?

- Retrieval models define the ranking (scores) of retrieved objects
- Several strands of models with many ways of estimating parameters stimulate IR research

Introduction Retrieval Models Related Models/Tasks Relationships between Retrieval Models Summary	Introduction Retrieval Models Related Models/Tasks Relationships between Retrieval Models Summary
Questions?	Brin, S. and Page, L. (1998). The anatomy of a large-scale hypertextual web search engine.
	In 7th International WWW Conference, Brisbane, Australia. Croft, W. and Harper, D. (1979). Using probabilistic models of document retrieval without relevance information. Journal of Documentation, 35:285–295. Robertson, S. and Sparck Jones, K. (1976). Relevance weighting of search terms. Journal of the American Society for Information Science, 27:129–146.
Thank you!	

17/18

IR4IP Tutorial IPI Confex, March 2009 Erik Graf and Thomas Roelleke User Interaction IPI Confex, March 2009

> Erik Graf University of Glasgow





1/19	2/19
	Introduction The Information Access Process Query Specification User Interfaces User Behavior Mining
Outline	Aim
2 The Information Access Process	
3 Query Specification	with a system
User Interfaces	
5 User Behavior Mining	

Introduction The Information Access Process Query Specification User Interfaces User Behavior Mining Introduction The Information Access Process Query Specification User Interfaces User Behavior Mining

## Models of Information Access

## Models of Information Access



#### The Problem

- Synonymity, Vocabulary mismatch (car, automobile, sedan, van)
- Polysemy (Chip, Java, etc ...)
- Vague and short user queries (2-3 terms average length for web queries)



Introduction The Information Access Process Query Specification User Interfaces User Behavior Mining

**Keyword Suggestion Techniques** 

Co-Occurrence analysis of terms
Pseudo-Relevance Feedback

Introduction The Information Access Process Query Specification User Interfaces User Behavior Mining

## **Query Specification**

Thesauri

Query Clustering

## Query Specification

Query Clustering

1	No.	Query Text	Clicked Documents
1	1.	law of thermodynamics	ID: 761571911 Title: Thermodynamics
			ID: 761571262 Title: Conservation Laws
2	2.	conservation laws	ID: 761571262 Title: Conservation Laws
			ID: 761571911 Title: Thermodynamics
3	3.	Newton law	ID: 761573959 Title: Newton, Sir Isaac
			ID: 761573872 Title: Ballistics
4	4.	Newton law	ID: 761556906 Title: Mechanics
			ID: 761556362 Title: Gravitation

	9/19		10/19
Introduction The Information Access Process Query Specification User Interfaces User Behavior Mining		Introduction The Information Access Process Query Specification User Interfaces User Behavior Mining	Information Item Representation
Query Transformation		Information item presentat	ion

#### Automatic query expansion

Users are generally reluctant to use suggested keywords. Therefore search engines often apply automated query expansion based on the following techniques:

- Thesauri, Co-occurrence analysis
- Pseudo-Relevance Feedback
- Query Clustering

- Summarization: Which snippet will allow for the best relevance judgment from the users side.
- Result page design: Tabbed result pages, amount of shown results.

Introduction The Information Access Process Query Specification User Interfaces User Behavior Mining

Information Item Representation

Introduction The Information Access Process Query Specification User Interfaces User Behavior Mining

Information Item Representation

### Information item presentation

### Automatically derived visualizations





#### Forms

- Implicit feedback: Automatically collected records of a users interaction with a system (e.g. 'click-through' data).
- Explicit feedback.

## Implicitly collected records

## Explicitly collected records

#### Applications

- Meta-data for the collection documents (a query-log field index)
- Collaborative filtering (a la Amazon: people who entered this query also clicked on ...)
- Query Clustering (Keyword suggestion)
- Personalization (Personal search history, topical preferences)

#### Applications

- Personalization
- Evaluation of retrieval performance.

18/19

17/19

Introduction The Information Access Process Query Specification User Interfaces User Behavior Mining

## Questions?

## Thank you!

## **IR4IP: Evaluation**



Fachgruppe Information Retrieval Gesellschaft für Informatik Herbstschule IR Schloss Dagstuhl 30. Sept. 2008

Thomas Mandl Informationswissenschaft Universität Hildesheim mandl@uni-hildesheim.de

## **IR4IP: Evaluation**

Slides by Thomas Mandl, IR Herbstschule 2008

**Updates by Thomas Roelleke** 

## Die Evaluierung von Information Retrieval Systemen



• Which System is better?

"There must be some fundamental understanding of what it means to be good and what it means to be better" (Bollmann/Cherniavsky 1983,3)

Mandl: Die Evaluierung von Information Retrieval Systemen



## **Recall und Precision**

 "The ability of the retrieval system to uncover relevant documents is known as the recall power of the system" (Lancaster 1968,55)

Recall =

Number of retrieved relevant documents Number of relevant documents

Precision =  $\frac{Number of retrieved relevant documents}{Number of retrieved documents}$ 

## **Role of Evaluation**

- Many components, models and optimisation techniques involved in a search system
- Effectiveness for a given (new) set of data difficult to forecast
- A general superiority of a single model or a single component is difficult to establish
- · Therefore, evaluate effectiveness
- · A holistic evaluation is difficult
- · Measure success/satisfaction of users?

Mandl: Die Evaluierung von Information Retrieval Systemen

## Example

Assume that for a given query, the following documents are relevant (10 relevant documents)

{d3, d5, d9, d25, d39, d44, d56, d71, d89, d123}

Suppose that the following documents are retrieved for that query:

rank	doc	precision	recall	rank	doc	precision	recall
1	d123	1/1	1/10	8	d129	1	$\leq t$
2	d84			9	d187		
3	d56	2/3	2/10	10	d25	4/10	4/10
4	d6			11	d48	$\sim$	
5	d8			12	d250	$\boldsymbol{\lambda}$	1L
6	d9	3/6	3/10	13	d113	//T	
7	d511			14	d3	5/14	5/10



## **Recall?**

- Users "feel" the precision
- Recall? Not "visible".
- Even with considerable effort difficult to determine precisely!
  - Number of relevant docs not know.
  - In particular problematic for queries where recall is important (e.g. crime investigations, legal applications, patent search)

Mandl: Die Evaluierung von Information Retrieval Systemen

## Prec at N

- Precision at N (10) documents
  - Clear interpretation
  - Reasonable for web retrieval
  - Little information about the system
  - Position of relevant documents not considered

## Which System is better?



## 2 Systems (A and B), 3 Topics

Mandl: Die Evaluierung von Information Retrieval Systemen

## "IR Psychology"

«The **unhappy customer**, on average, **will tell 27** other people ...»

 $\rightarrow$  Bad news travels fast.

Site search needs to be robust Avoid bad outliers!

for as many queries as possible

for as many measures as possible

Mandl: Die Evaluierung von Information Retrieval Systemen Credit to Jacques Savoy for slide

## **NTCIR**



- Cross-lingual IR asian languages
- Tokio
  - National Institute for Informatics
- Tasks
  - Cross-lingual
    - Chinesisch, Japanisch, Koreanisch -> Englisch
  - Patent-Retrieval
  - Web-Retrieval
  - Question Answering

