# Semistructured Data Search

Krisztian Balog University of Stavanger

Promise Winter School 2013 | Bressanone, Italy, February 2013





#### - Semistructured data

- Lack of fixed, rigid schema
- No separation between the data and the schema, self-describing structure (tags or other markers)

## Motivation

- Supporting users who cannot express their need in structured query languages
  - SQL, SPARQL, Inquery, etc.
- Dealing with heterogeneity
  - Users are unaware of the schema of the data
  - No single schema to the data

## Semistructured data

#### - Advantages

- The data is not constrained by a fixed schema
- Flexible (the schema can easily be changed)
- Portable
- Possible to view structured data as semistructured

#### - Disadvantages

- Queries are less efficient than in a constrained structure

## In this talk

- How to exploit the structure available in the data for retrieval purposes?
- Different types of structure
  - Document, query, context
- Working in a Language Modeling setting
- Number of different tasks
  - Retrieving entire documents
    - I.e., no element-level retrieval
  - Textual document representation is readily available
    - No aggregation over multiple documents/sources

## **Incorporating structure**



# Preliminaries

Language modeling

## Language Modeling

- Rank documents *d* according to their likelihood of being relevant given a query q: P(d|q)





## Language Modeling



## Example

In the town where I was born, Lived a man who sailed to sea, And he told us of his life, In the land of submarines,

So we sailed on to the sun, Till we found the sea green, And we lived beneath the waves, In our yellow submarine,

We all live in yellow submarine, yellow submarine, yellow submarine, We all live in yellow submarine, yellow submarine, yellow submarine.



## **Empirical document LM**

 $P(t|d) = \frac{n(t,d)}{|d|}$ 



## Alternatively...



## Scoring a query

 $q = \{\text{sea}, \text{submarine}\}$ 

 $P(q|d) = P(\text{``sea''}|\theta_d) \cdot P(\text{``submarine''}|\theta_d)$ 

## Scoring a query

 $q = \{\text{sea}, \text{submarine}\}$ 

 $P(q|d) = \underbrace{P(\text{``sea''}|\theta_d)}_{\substack{0.03602}} \cdot P(\text{``submarine''}|\theta_d)$   $\underbrace{P(q|d) = \underbrace{P(\text{``sea''}|\theta_d)}_{\substack{0.1 \\ (1 - \lambda)P(\text{``sea''}|d) + \lambda P(\text{``sea''}|C)}}_{\substack{0.0002 \\ (1 - \lambda)P(\text{``sea''}|d) + \lambda P(\text{``sea''}|C)}}$ 

t	P(t d)
submarine	0.14
sea	0.04

t	P(t C)
submarine	0.0001
sea	0.0002
•••	

## Scoring a query

 $q = \{\text{sea}, \text{submarine}\}$ 

0.04538 0.03602 0.12601  $P(q|d) = P(\text{``sea''}|\theta_d) \cdot \underbrace{P(\text{``submarine''}|\theta_d)}_{0.14} \quad \underbrace{P(\text{``submarine''}|\theta_d)}_{0.14} \quad \underbrace{P(\text{``submarine''}|\theta_d)}_{0.10001} \quad \underbrace{P(\text{`$ 

t	P(t d)
submarine	0.14
sea	0.04

t	P(t C)
submarine	0.0001
sea	0.0002
•••	

# Part I

**Document structure** 

## In this part

- Incorporate document structure into the document language model
  - Represented as *document fields*

$$\begin{split} P(d|q) \propto P(q|d)P(d) = P(d) \prod_{t \in q} \underbrace{P(t|\theta_d)}_{t \in q}^{n(t,q)} \end{split} \\ \end{split}$$

#### **Use case** Web document retrieval

promise Participative Research labOratory and Multilingual Information System	SE for Multimedia ms Evaluation	🎤 Log-in
Overview Achievements Use cases Publication	ons Events CLEF Media Center Contacts	Q Search
Events Winter School 2013		
PROMISE Winter School 2013 Bressanone, Italy		
Winter School 2013	PROMISE Winter	School 2013
Programme     Lecturers	Bridging between Information F	Retrieval and Databases
<ul> <li>Venue</li> <li>Registration and Accomodation</li> <li>Sponsor and Patronage</li> <li>Flyer</li> </ul>	Bressanone, Italy 4 - 8	February 2013
Important Dates Registration Deadline (extended): 28 <sup>th</sup>	The aim of the PROMISE Winter School 2013 on "Bridging between participants a grounding in the core topics that constitute the multi- unstructured, semistructured, and structured information. The scho- from invited speakers who are recognized experts in the field. The or senior researchers such as post-doctoral researchers form the fields.	en Information Retrieval and Databases" is to give disciplinary area of information access and retrieval to bol is a week-long event consisting of guest lectures school is intended for PhD students, Masters students fields of databases, information retrieval, and related

### Web document retrieval Unstructured representation

PROMISE Winter School 2013 Bridging between Information Retrieval and Databases Bressanone, Italy 4 - 8 February 2013 The aim of the PROMISE Winter School 2013 on "Bridging between Information Retrieval and Databases" is to give participants a grounding in the core topics that constitute the multidisciplinary area of information access and retrieval to unstructured, semistructured, and structured information. The school is a week-long event consisting of guest lectures from invited speakers who are recognized experts in the field. The school is intended for PhD students, Masters students or senior researchers such as post-doctoral researchers form the fields of databases, information retrieval, and related fields.

[...]

### Web document retrieval HTML source

<html>

<head>

```
<title>Winter School 2013</title>
```

```
<meta name="keywords" content="PROMISE, school, PhD, IR, DB, [...]" />
<meta name="description" content="PROMISE Winter School 2013, [...]" />
</head>
```

<body>

```
<h1>PROMISE Winter School 2013</h1>
```

```
<h2>Bridging between Information Retrieval and Databases</h2>
<h3>Bressanone, Italy 4 - 8 February 2013</h3>
The aim of the PROMISE Winter School 2013 on "Bridging between
Information Retrieval and Databases" is to give participants a grounding
in the core topics that constitute the multidisciplinary area of
information access and retrieval to unstructured, semistructured, and
structured information. The school is a week-long event consisting of
guest lectures from invited speakers who are recognized experts in the
field. The school is intended for PhD students, Masters students or
senior researchers such as post-doctoral researchers form the fields of
databases, information retrieval, and related fields.
```

</body>

</html>

#### **Web document retrieval** Fielded representation based on HTML markup

- title: Winter School 2013
- meta: PROMISE, school, PhD, IR, DB, [...]
  PROMISE Winter School 2013, [...]
- headings: PROMISE Winter School 2013
  Bridging between Information Retrieval and Databases
  Bressanone, Italy 4 8 February 2013
- **body:** The aim of the PROMISE Winter School 2013 on "Bridging between Information Retrieval and Databases" is to give participants a grounding in the core topics that constitute the multidisciplinary area of information access and retrieval to unstructured, semistructured, and structured information. The school is a weeklong event consisting of guest lectures from invited speakers who are recognized experts in the field. The school is intended for PhD students, Masters students or senior researchers such as postdoctoral researchers form the fields of databases, information retrieval, and related fields.

## **Fielded Language Models** [Ogilvie & Callan, SIGIR'03]

- Build a separate language model for each field
- Take a linear combination of them



## Field Language Model



## Fielded Language Models Parameter estimation

- Smoothing parameter
  - Dirichlet smoothing with avg. representation length
- Field weights
  - Heuristically (e.g., proportional to the length of text content in that field)
  - Empirically (using training queries)
    - Computationally intractable for more than a few fields

## Example

 $\begin{aligned} q &= \{\text{IR, winter, school}\} \\ \text{fields} &= \{\text{title, meta, headings, body}\} \\ \mu &= \{0.2, 0.1, 0.2, 0.5\} \end{aligned}$ 

$$P(q|\theta_d) = \underbrace{P(``IR" | \theta_d)}_{\downarrow} \cdot P(``winter" | \theta_d) \cdot P(``school" | \theta_d)$$

$$P(``IR" | \theta_d) = 0.2 \cdot P(``IR" | \theta_{d_{title}})$$

$$+ 0.1 \cdot P(``IR" | \theta_{d_{meta}})$$

$$+ 0.2 \cdot P(``IR" | \theta_{d_{headings}})$$

$$+ 0.2 \cdot P(``IR" | \theta_{d_{headings}})$$

### **Use case Entity retrieval in RDF data**

#### Create account 🔒 Log in

Q

WikipediA The Free Encyclopedia

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 العريبة Беларуская

Article Talk

#### Audi A4

From Wikipedia, the free encyclopedia

The Audi A4 is a line of compact executive cars produced since late 1994 by the German car manufacturer Audi, a subsidiary of the Volkswagen Group.

The A4 has been built in four generations and is based on Volkswagen's B platform. The first generation A4 succeeded the Audi 80. The automaker's internal numbering treats the A4 as a continuation of the Audi 80 lineage, with the initial A4 designated as the B5-series, followed by the B6, B7, and the current B8. The B8 A4 is built on the Volkswagen Group MLB platform shared with many other Audi models and potentially one Porsche model within Volkswagen Group.<sup>[2]</sup>

The Audi A4 automobile layout consists of a longitudinally oriented engine at the front, with transaxle-type transmissions mounted at the rear of the engine. The cars are front-wheel drive, or on some models, "guattro" all-wheel drive.

The A4 is available as a saloon/sedan and estate/wagon. The second (B6) and third generations (B7) of the A4 also had a convertible version, but the B8 version of the convertible became a variant of the Audi A5 instead as Audi got back into the compact executive coupé segment. The facebook fans of the Audi A4 page are more than 870,000.

Contents [show]

Edit View history Read

Search

Audi A4



Manufacturer	Audi
Production	1994-present
Assembly	Ingolstadt, Germany Changchun, China <sup>[1]</sup> Tokyo, Japan (AMA; B5 only) Jakarta, Indonesia (Garuda Mataram Motor; B5 & B8) Solomonovo, Ukraine (Eurocar; B7 only) Aurangabad, India
Dradaaaaaa	Audi 00

#### Predecessor Audi 80

Class	Compact executive car (globally)
	the set of the set of the set of the set of the

#### **Use case** Entity retrieval in RDF data

#### dbpedia:Audi\_A4

foaf:name rdfs:label	Audi A4 Audi A4
rdfs:comment	The Audi A4 is a compact executive car produced since late 1994 by the German car manufacturer Audi, a subsidiary of the
dbpprop:production	Volkswagen Group. The A4 has been built [] 1994 2001 2005 2008
rdf:type	dbpedia-owl:MeanOfTransportation dbpedia-owl:Automobile
dbpedia-owl:manufacturer	dbpedia:Audi
dbpedia-owl:class	dbpedia:Compact_executive_car
owl:sameAs	freebase:Audi A4
is <b>dbpedia-owl:predecessor</b> of	dbpedia:Audi_A5
is <b>dbpprop:similar</b> of	dbpedia:Cadillac_BLS

## **Hierarchical Entity Model** [Neumayer et al., ECIR'12]

- Number of possible fields is huge
  - It is not possible to optimise their weights directly
- Entities are sparse w.r.t. different fields
  - Most entities have only a handful of predicates
- Organise fields into a 2-level hierarchy
  - Field types (4) on the top level
  - Individual fields of that type on the bottom level
- Estimate field weights
  - Using training data for field types
  - Using heuristics for bottom-level types

## **Two-level hierarchy**

Name	foaf:name	Audi A4
	rdfs:label	Audi A4
Attributes { Out-relations {	rdfs:comment	The Audi A4 is a compact executive car
		produced since late 1994 by the German car
		manufacturer Audı, a subsıdıary of the
		Volkswagen Group. The A4 has been built []
	dbpprop:production	1994
		2001
		2005
		2008
	rdf:type	dbpedia-owl:MeanOfTransportation
		dbpedia-owl:Automobile
	dbpedia-owl:manufacturer	dbpedia:Audi
	dbpedia-owl:class	dbpedia:Compact_executive_car
	owl:sameAs	freebase:Audi A4
In-relations -	is <b>dbpedia-owl:predecessor</b> of	dbpedia:Audi_A5
	is <b>dbpprop:similar</b> of	dbpedia:Cadillac_BLS

## **Hierarchical Entity Model** [Neumayer et al., ECIR'12]



# Field generation

#### - Uniform

- All fields of the same type are equally important
- Length
  - Proportional to field length (on the entity level)
- Average length
  - Proportional to field length (on the collection level)
- Popularity
  - Number of documents that have the given field

## **Comparison of models**







Unstructured document model

Fielded document model

Hierarchical document model

#### **Use case** Finding movies in IMDB data



### **Use case** Finding movies in IMDB data

```
<title>The Transporter</title>
<year>2002</year>
<language>English</language>
<genre>Action</genre>
<genre>Crime</genre>
<genre>Thriller</genre>
<country>USA</country>
<actors>
   <actor>Jason Statham</actor>
   <actor>Matt Schulze</actor>
   <actor>François Berléand</actor>
   <actor>Ric Young</actor>
   <actress>0i Shu/actress>
</actors>
<team>
   <director>Louis Leterrier</director>
   <director>Corey Yuen</director>
   <writer>Luc Besson</writer>
   <writer>Robert Mark Kamen</writer>
   producer>Luc Besson</producer>
   <cinematographer>Pierre Morel</cinematographer>
</team>
```

## **Probabilistic Retrieval Model for Semistructured data** [Kim et al., ECIR'09]

- Find which document field each query term may be associated with
- Extending [Ogilvie & Callan, SIGIR'03]


# **PRMS**Mapping probability



#### Prior field probability

Probability of mapping the query term to this field before observing collection statistics

# **PRMS**Mapping example



# Part 2

Query structure

# Structured query representations

- Query may have a semistructured representation, i.e., multiple fields
- Examples
  - TREC Genomics track
    - (1) gene name, (2) set of symbols
  - TREC Enterprise track, document search task
    - (1) keyword query, (2) example documents
  - INEX Entity track
    - (1) keyword query, (2) target categories, (3) example entities
  - TREC Entity track
    - (1) keyword query, (2) input entity, (3) target type

#### Use case

#### **Enterprise document search (TREC 2007)**

- Task: create an overview page on a given topic
  - Find documents that discuss the topic in detail

#### cancer risk



```
<topic>
<num>CE-012</num>
<query>cancer risk</query>
<narr>
Focus on genome damage and therefore cancer risk in humans.
</narr>
<page>CSIR0145-10349105</page>
<page>CSIR0140-15970492</page>
<page>CSIR0139-07037024</page>
<page>CSIR0138-00801380</page>
</topic>
```



# Query modeling

- Aims
  - Expand the original query with additional terms
  - Assign the probability mass non-uniformly



### **Retrieval model**

- Maximizing the query log-likelihood provides the same ranking as minimising KL-divergence
  - Assuming uniform document priors



# Estimating the query model

- Baseline maximum-likelihood

$$P(t|\theta_q) = P(t|q) = \frac{n(t,q)}{|q|}$$

- Query expansion using relevance models [Lavrenko & Croft, SIGIR'01]

$$P(t|\theta_q) = (1 - \lambda) P(t|\hat{q}) + \lambda P(t|q)$$

$$\textbf{Expanded query model}$$
Based on term co-occurrence statistics
$$P(t|\hat{q}) \approx \frac{P(t, q_1, \dots, q_k)}{\sum_{t'} P(t', q_1, \dots, q_k)}$$

#### **Sampling from examples** [Balog et al., SIGIR'08]



#### **Sampling from examples** Importance of a sample document

- Uniform P(d|S) = 1/|S|
  - All sample document are equally important
- Query-biased  $P(d|S) \propto P(d|q)$ 
  - Proportional to the document's relevance to the (original) query
- Inverse query-biased  $P(d|S) \propto 1 P(d|q)$ 
  - Reward documents that bring in new aspects (not covered by the original query)

#### **Sampling from examples** Estimating term importance

- Maximum-likelihood estimate

$$P(t|d) = \frac{n(t,d)}{|d|}$$

- Smoothed estimate

 $P(t|d) = P(t|\theta_d) = (1 - \lambda)P(t|d) + \lambda P(t|C)$ 

- Ranking function by [Ponte, 2000]

$$P(t|d) = \frac{s(t)}{\sum_{t'} s(t')} \quad s(t) = \log \frac{P(t|d)}{P(t|C)}$$

### Use case

#### **Entity retrieval in Wikipedia (INEX 2007-09)**

- Given a query, return a ranked list of entities
  - Entities are represented by their Wikipedia page
- Entity search
  - Topic definition includes target categories
- List search
  - Topic definition includes example entities

#### Titanic (1997 film)

From Wikipedia, the free encyclopedia

*Titanic* is a 1997 American epic romance and disaster film directed, written, co-produced, and co-edited by James Cameron. A fictionalized account of the sinking of the RMS *Titanic*, it stars Leonardo DiCaprio as Jack Dawson and Kate Winslet as Rose DeWitt Bukater, members of different social classes who fall in love aboard the ship during its ill-fated maiden voyage. Although the central roles and love story are fictitious, some characters are based on genuine historical figures. Gloria Stuart portrays the elderly Rose, who narrates the film in a modern-day framing device, and Billy Zane plays Cal Hockley, the overbearing fiancé of the younger Rose. Cameron saw the love story as a way to engage the audience with the real-life tragedy.

Production on the film began in 1995, when Cameron shot footage of the actual *Titanic* wreck. The modern scenes were shot on board the *Akademik Mstislav Keldysh*, which Cameron had used as a base when filming the actual wreck. A reconstruction of the *Titanic* was built at Playas de Rosarito, Baja California, and scale models and computer-generated imagery were also used to recreate the sinking. The film was partially funded by Paramount Pictures and 20th Century Fox – respectively, its American and international distributor – and at the time, it was the most expensive film ever made, with an estimated budget of US\$200 million.<sup>[3][4][5][6]</sup>

The film was originally scheduled to open on July 2, 1997, however, post-production delays pushed back its release to December 19 instead.<sup>[7]</sup> *Titanic* was an enormous critical and commercial success. It was nominated for fourteen Academy Awards, eventually winning eleven, including Best Picture and Best Director.<sup>[8]</sup> It became the highest-grossing film of all time, with a worldwide gross of over \$1.8 billion, and remained so for twelve years until Cameron's next directorial effort, *Avatar*, surpassed it in 2010.<sup>[9][10]</sup> *Titanic* also has been ranked as the sixth best epic film of all time in AFI's 10 Top 10 by the American Film Institute.<sup>[11]</sup> The film is due for theatrical re-release in 2012 after Cameron completes its conversion into 3-D.<sup>[12]</sup>



Theotrical poster

Categories: 1997 films I American films I English-language films I American disaster films I Best Drama Picture Golden Globe winners I Best Picture Academy Award winners I Best Song Academy Award winners I Films directed by James Cameron I Films set in 1912 I Films that won the Best Sound Mixing Academy Award I Films that won the Best Visual Effects Academy Award I Films whose art director won the Best Art Direction Academy Award I Films whose cinematographer won the Best Cinematography Academy Award I Films whose director won the Best Director Academy Award I Films whose director Golden Globe I Films whose editor won the Best Film Editing Academy Award I Epic films I RMS Titanic I Romantic epic films I Romantic period films I Seafaring films based on actual events I Films shot in Nova Scotia I Films shot in Vancouver I Paramount films I 20th Century Fox films I Lightstorm Entertainment films I 2-D films converted to 3-D

## **Example query**

<title>Movies with eight or more Academy Awards</title> <categories>

<category id="45168">best picture oscar</category>
 <category id="14316">british films</category>
 <category id="2534">american films</category>
</categories>

#### Using categories for retrieval

- As a separate document field
- Filtering
- Similarity between target and entity categories
  - Set similarity metrics
  - Content-based (concatenating category contents)
  - Lexical similarity of category names

## Also related to categories

#### - Category expansion

- Based on category structure
- Using lexical similarity of category names
- Ontology-base expansion
- Generalisation
  - Automatic category assignment





# Part 3

Contextual structures

# Other types of structure

- Link structure
- Linguistic structure
- Social structures

### Link structure

Source: Wikipedia

### Link structure

- Aim
  - High number of inlinks from pages relevant to the topic and not many incoming links from other pages
- Retrieve an initial set of documents, then rerank



#### **Document link degree** [Kamps & Koolen, ECIR'08]



#### **Relevance propagation** [Tsikrika et al., INEX'06]

- Model a user (random surfer) that after seeing the initial set of results
  - Selects one document and reads its description
  - Follows links connecting entities and reads the descriptions of related entities
  - Repeats it N times

$$P_{0}(d) = P(q|d)$$

$$P_{i}(d) = \underbrace{P(q|d)}_{\downarrow} P_{i-1}(d) + \underbrace{\sum_{d' \to d}}_{\downarrow} (1 - P(q|d')) \underbrace{P(d'|d)}_{\downarrow} P_{i-1}(d')$$
The probability of staying at the node equals to its relevance to the query
$$Outgoing links from d' to d$$
Transition probabilities

#### **Relevance propagation** [Tsikrika et al., INEX'06]

- Weighted sum of probabilities at different steps

$$P(d) \propto \mu_0 P_0(d) + (1 - \mu_0) \sum_{i=1}^N \mu_i P_i(d)$$



# **Translation Model**

[Berger & Lafferty, SIGIR'99]

$$P(q|d) = \prod_{t \in q} \left( \sum_{w \in V} \underbrace{P(t|w)}_{w \in V} P(w|\theta_d) \right)^{n(t,q)}$$
  
**Translation model**  
Probability that word w can

"semantically translated" to word qi

- Obtaining the translation model
  - Exploiting WordNet word co-occurrences [Cao et al., SIGIR'05]

#### **Use case** Expert finding

- Given a keyword query, return a ranked list of people who are experts on the given topic
- Content-based methods can return the most knowledgeable persons
- However, when it comes to contacting an expert, social and physical proximity matters

### **Expert profile**

#### <person>

```
<anr>710326</anr>
<name>Toine M. Bogers</name>
<name>Toine Bogers</name>
<name>A. M. Bogers</name>
<job>PhD student</job>
<faculty>Faculty of Humanities</faculty>
<department>Department of Communication and Information Sciences</department>
<room>Room D 348</room>
<address>P.O. Box 90153, NL-5000 LE Tilburg, The Netherlands</address>
<tel>+31 13 466 245</tel>
<fax>+31 13 466 289</fax>
<homepage>http://ilk.uvt.nl/~toine</homepage>
<cemail>A.M.Bogers@uvt.nl</cemail>
```

```
<publications>
  <publication arno-id="1212">
     <title>Design and Implementation of a University-wide
        Expert Search Engine</title>
        <author>R. Liebregts and T. Bogers</author>
        <year>2009</year>
        <booktitle>Proceedings of the 31st European Conference on Information
Retrieval</booktitle>
        </publication>
        [...]
        </publications>
```

```
[...] </person>
```

#### **User-oriented model for EF** [Smirnova & Balog, ECIR'11]

$$S(e|u,q) = (1-\lambda)\underbrace{K(e|u,q)}_{\downarrow} + \lambda \underbrace{T(e|u)}_{\downarrow}$$

Knowledge gain

Difference between the knowledge of the expert and that of the user on the query topic

**Contact time** 

Distance between the user and the expert in a social graph

## **Social structures**

#### **Organisational hierarchy**

```
<person>
  <anr>710326</anr>
  <name>Toine M. Bogers</name>
  [...]
  <faculty>Faculty of Humanities</faculty>
  <department>Department of Communication
    and Information Sciences</department>
  [...]
  </person>
```



## **Social structures**

#### **Geographical information**





#### **Social structures** Co-authorship

```
<person>
  <anr>710326</anr>
 <name>Toine M. Bogers</name>
  [...]
  <publications>
    <publication arno-id="1212">
      <title>
        Design and Implementation of a
        University-wide Expert Search Engine
      </title>
      <author>
        R. Liebregts and T. Bogers
      </author>
      <year>2009</year>
      <booktitle>
        Proceedings of the 31st European
        Conference on Information Retrieval
      </booktitle>
    </publication>
    [...]
  </publications>
  [...]
</person>
```



### **Social structures**



# Summary

- Different types of structure
  - Document structure
  - Query structure
  - Contextual structures
- Probabilistic IR and statistical Language Models yield a principled framework for representing and exploiting these structures

# **Questions?**

Contact | @krisztianbalog | krisztianbalog.com