

Indexing of Political Texts: Combining Different Sources of Evidence

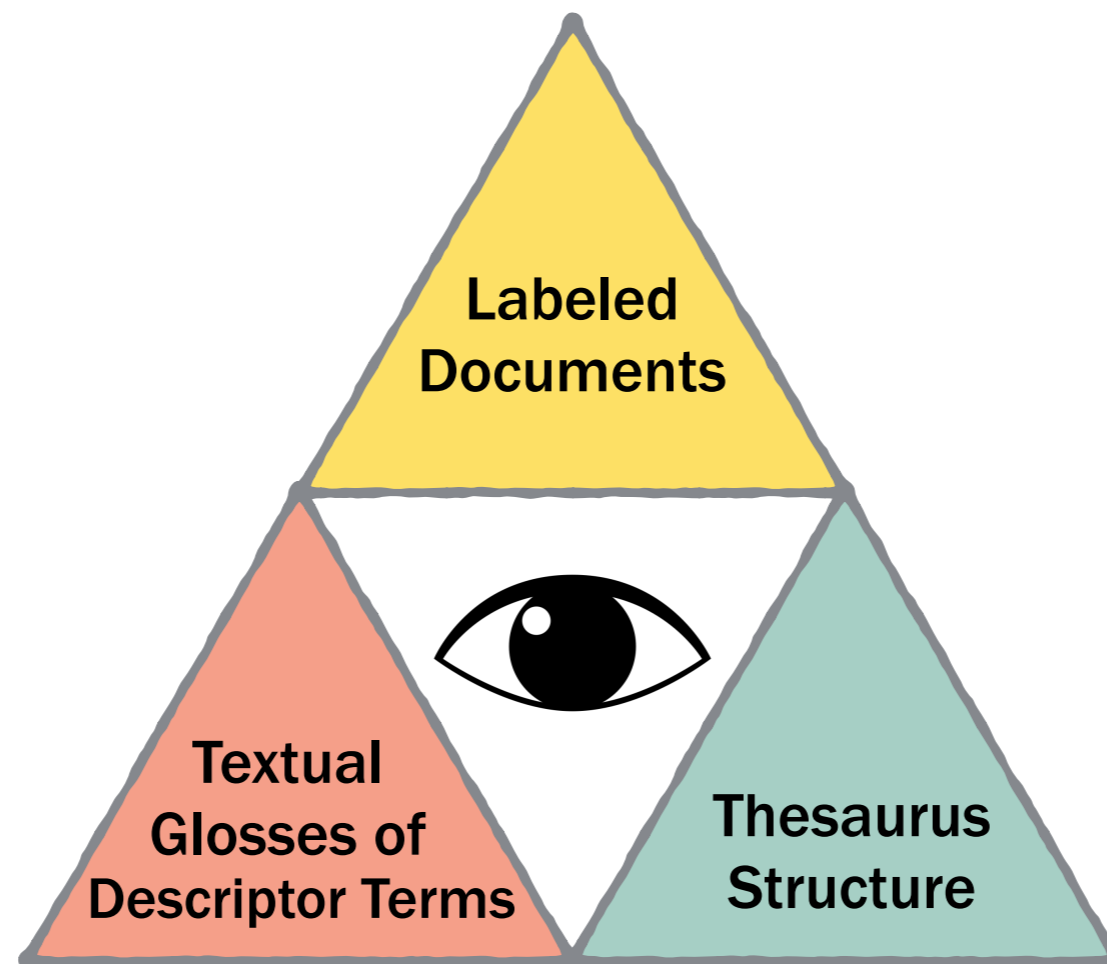
Mostafa Dehghani, Hosein Azarbonyad, Maarten Marx, Jaap Kamps

Research Questions

- What are the different possible **sources of evidence** for automatic indexing of political text?
- How effective is a **learning to rank** (LTR) approach integrating a variety of sources of information as features?
- What is the **relative importance** of each of these sources of information for indexing political text?
- Can we build a **lean-and-mean system** that approximate the effectiveness of the large LTR system?

Sources of Evidence

- There are different sources of evidence for the selection of appropriate indexing terms for political documents:



Sources of Evidence Combination

- **Learning to rank**

- Mapping the problem to an information retrieval problem

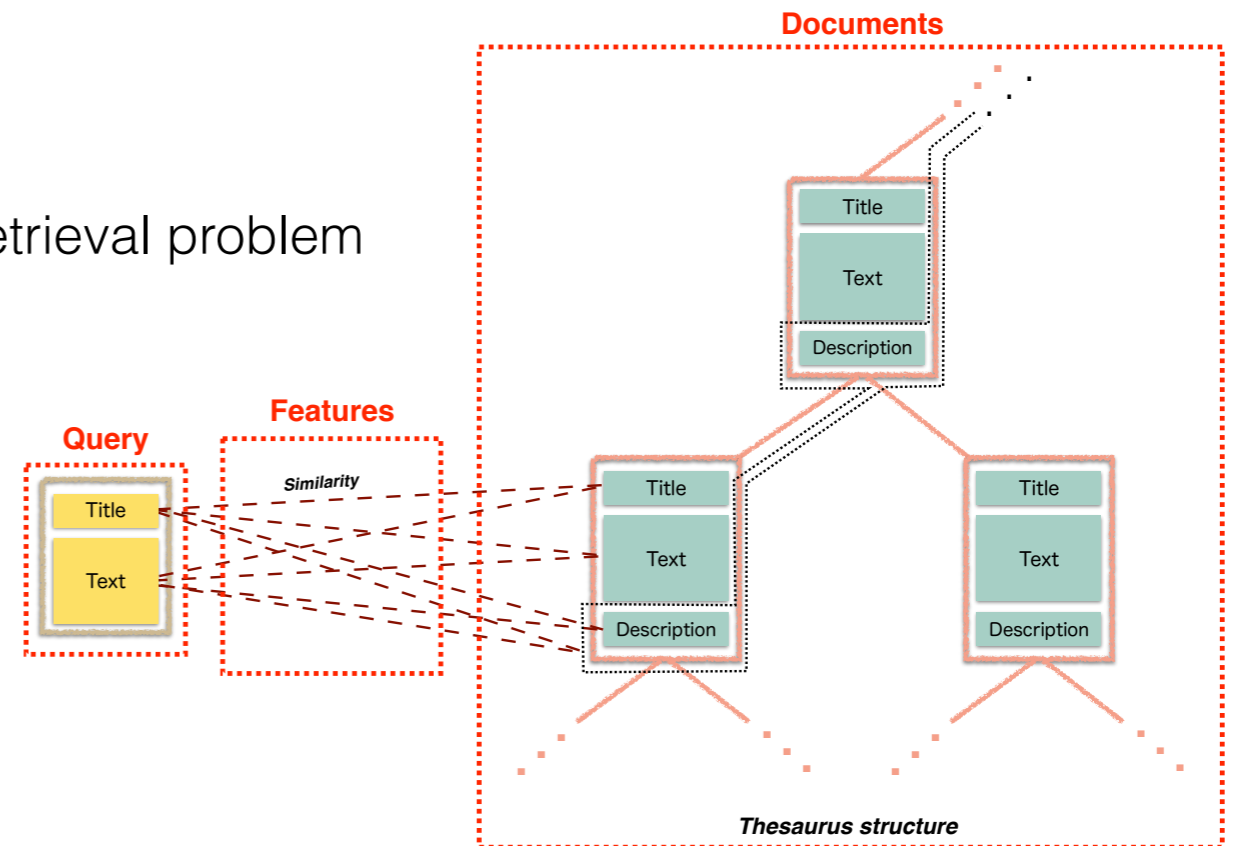
- Documents as *queries*

- Descriptor terms as *document*

- Features:

- Query dependent features (Similarity of query and documents)

- Query independent features (Prior knowledge)



Features

- **Query dependent features**

- Model of documents:

- **$Model_D = \langle M(title_D), M(text_D) \rangle$**

- Model of descriptor terms:

- **$Model_DT = \langle M(title_DT), M(text_DT), M(gloss_DT), M(anc-gloss_DT) \rangle$**

- Eight possible combinations of a document and descriptor term submodel

- Different IR Similarity Measures (Language modelling similarity based on KL-divergence using Dirichlet smoothing, the same run using Jelinek-Mercer smoothing, Okapi-BM25)

- **Query independent features**

- Popularity
- Generality
- Ambiguity

Experiments

- Data Collection:

- JRC-Acquis dataset



- The documents of this corpus have been manually labeled with *EuroVoc* concept



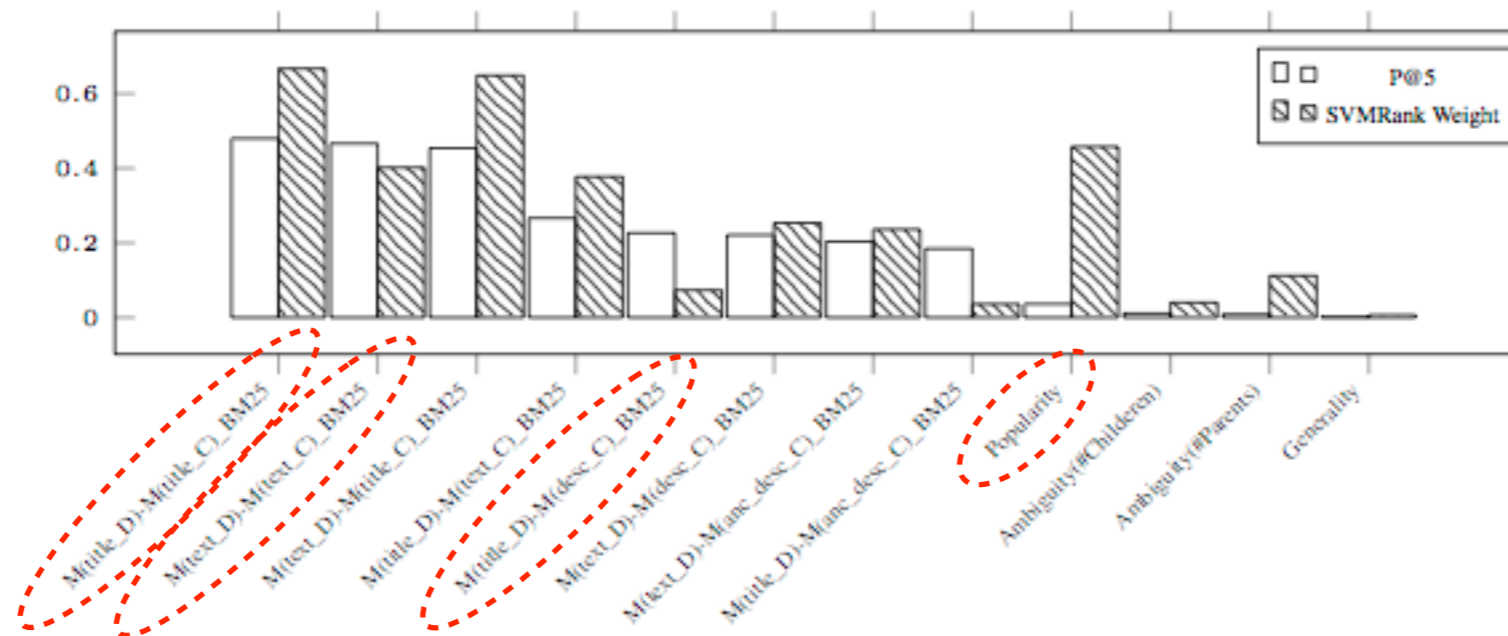
- **Effectiveness of LTR:**

- Performance of LTR over JEX as baseline system

Method	P@5 (%Diff.)	Recall@5 (%Diff.)
JEX	0.4353	0.4863
BM25-TITLES	0.4798 (10%) [^]	0.5064 (4%) [^]
LTR-ALL	0.5206 (20%) [^]	0.5467 (12%) [^]

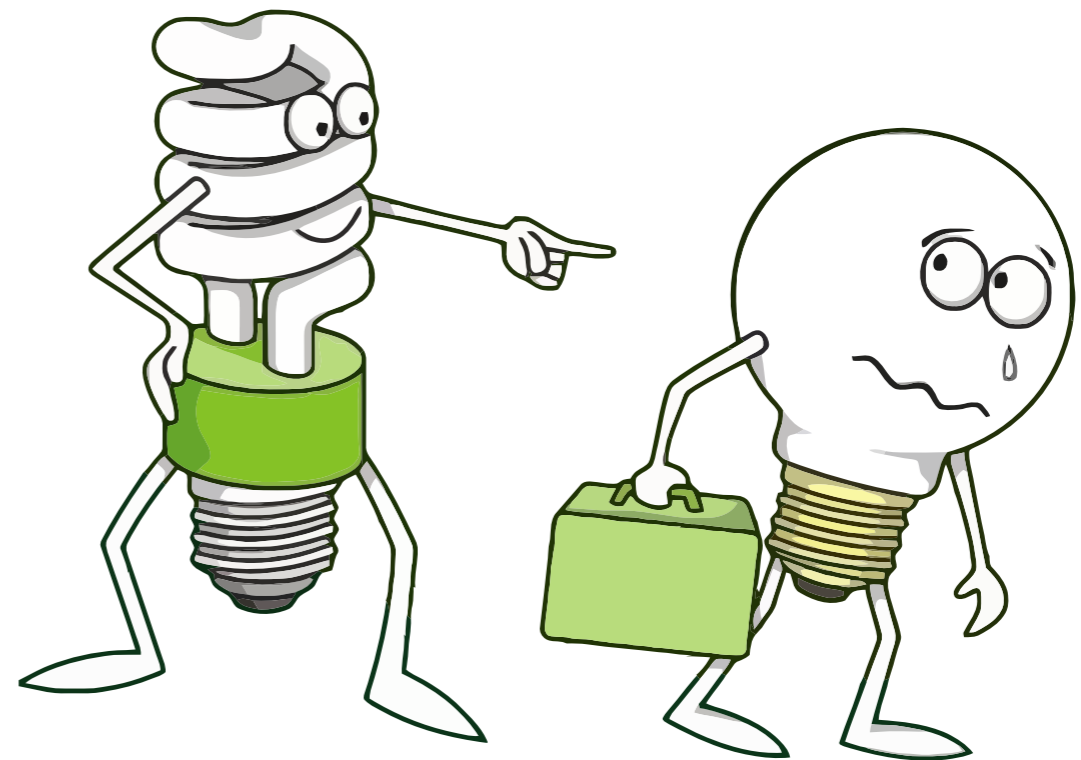
Experiments

- Feature Analysis: **Importance of Different Information Sources**



Experiments

- **Lean-and-Mean** system (Only 4 features):
 - *text submodel of documents with **all text***
 - ***titles** only*
 - *textual **glosses** of descriptor terms*
 - ***popularity** of descriptor terms*



Method	P@5 (%Diff.)	Recall@5 (%Diff.)
LTR-ALL	0.5206 (-)	0.5467 (-)
LTR-TTGP	0.5058 (-3%)	0.5301 (-3%)

Conclusion

- Our main findings:
 - Using a learning to rank (**LTR**) approach **integrating all features** significantly improves the performance of indexing than previous systems.
 - Analysis of feature weights:
 - **relative importance of various sources of evidence**
- **Lean-and-mean system:**
 - performs at **97%** of the large LTR model

Thank you