

Latent Topic Feedback for Information Retrieval

David Andrzejewski and David Buttler



Summary

We propose a method for improving ad-hoc information retrieval by allowing explicit user feedback over topics automatically learned from the corpus using the Latent Dirichlet Allocation (LDA) [1] model. This capability may be especially useful within organizations with specialized domains or limited resources. Experiments on TREC data with simulated user feedback show improved retrieval performance, in addition to the informational benefits of the displayed topics.

Information retrieval in challenging environments

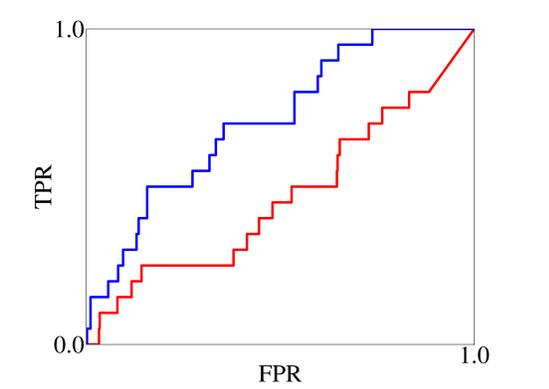
Conditions	Impaired techniques
Non-expert user	keyword queries
Lack of metadata	faceted search
Specialized domain	WordNet
Small user base	query log mining
	relevance feedback

Experimental results

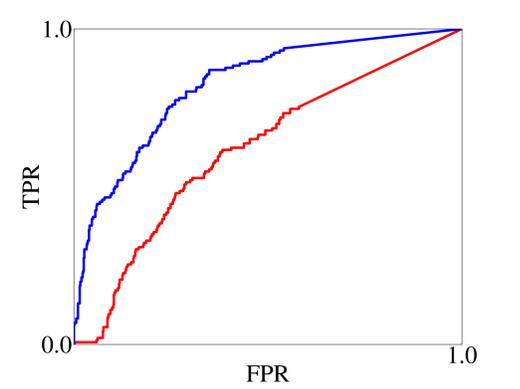
Queries and topics

ROC curves
BASELINE
VS
FEEDBACK

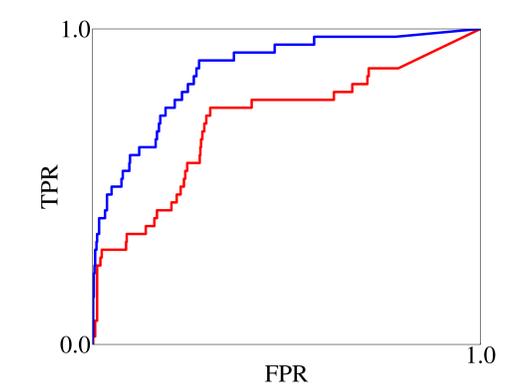
“euro opposition”
(**Emu**) economic monetary union
Maastricht treaty, member states
European, Europe, Community, Emu



“King Hussein, peace”
(**Amman**) Majesty King Husayn
al Aqabah, peace process
Jordan, Jordanian, Amman, Arab

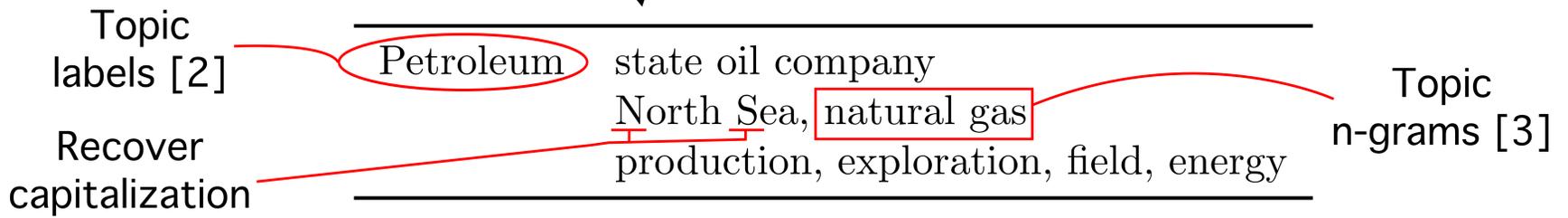


“law enforcement dogs”
(**heroin**) seized kg cocaine
drug traffickers, kg heroin
police, arrested, drugs, marijuana



1) Topic representation

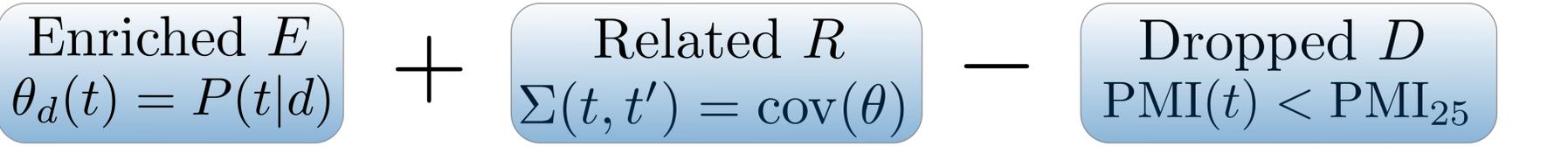
Typical "Top N" representations of learned topics can be difficult for users to interpret. We combine several methods to construct a more easily understood topic summary.



2) Topic selection

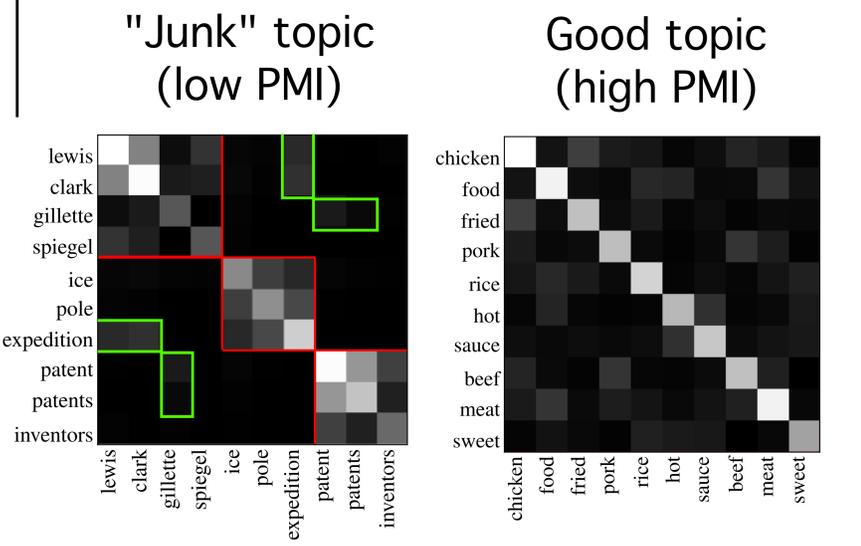
Topics enriched in the top documents returned by the original query

We may have hundreds of learned topics. We use the following procedure to select a handful of topics to present to the user as feedback candidates.



3) Incorporating feedback

When a user marks a topic as relevant, we expand the original query using the most probable words from that topic. In order to preserve query intent, the original query is given more weight.



References

- [1] D. Blei, A. Ng, and M. Jordan. Latent Dirichlet Allocation. (JMLR 2003)
- [2] J.H. Lau, D. Newman, S. Karimi, and T. Baldwin. Best topic word selection for topic labelling. (COLING 2010)
- [3] D. Blei and J. Lafferty. Visualizing topics with multi-word expressions. Technical report, 2009. arXiv:0907.1013v1 [stat.ML]
- [4] D. Newman, Y. Noh, E. Talley, S. Karimi, and T. Baldwin. Evaluating topic models for digital libraries. (JCDL 2010)

Acknowledgments

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344 (LLNL-CONF-471258). LLNL-POST-495651