

Semantic-Powered Research Profiling

Zhixiong Zhang
National Science Library
Chinese Academy of Sciences
Beijing, China
+86-10-82629426
zhangzhx@mail.las.ac.cn

Ying Ding
Indiana University
1320 E 10th
Bloomington, IN
+1-812-855-5388
dingying@indiana.edu

Na Hong
National Science Library
Chinese Academy of Sciences
Beijing, China
+86-10-82629426
hongn@mail.las.ac.cn

ABSTRACT

Research profiling is a widely-adopted method to monitor research development and rank research performance. This paper describes a novel infrastructure to generate semantic-powered research profiling for research fields, organizations and individuals. It crawls related websites and news feeds, extracts research terms, research objects and relations from them and uses the proposed Research Ontology to model them into RDF triples to facilitate semantic queries and semantic mining on burst detection, hot topic detection, dynamics of research, and relation mining. The authors implement a research profiling experiment in Artificial Intelligence area to show the effectiveness of the research profiling based on semantic mining.

Categories and Subject Descriptors

I.2.4 [Knowledge Representation Formalisms and Methods]: Semantic Networks

General Terms

Measurement, Design

Keywords

Research Profiling, Semantic mining, Knowledge component Extraction, Visualization

1. INTRODUCTION

Research profiling is a widely-adopted method to monitor research development and rank research performance within a certain research field [1] [2]. It gathers related research materials via automatic crawling on organizational websites, news feeds, personal websites, online journals and related databases, extracts valuable data by automatic information extraction tools, and creates evaluation metrics based on co-occurrence analysis and other research policy indicators. Through the daily monitoring the selected research field, it can obtain a “big picture” on the research activity, understand the research community, gain insight into how innovation is progressing, and map (graphically represent) topical interrelationships for a whole research field.

Related work on research profiling include [1],[2],and [3]. Most of them based on traditional view and applied to structured text literature. Although in unstructured data content type aspect, some attempts have been made, however, it focused on shallow statistic analysis. For example, Alan Porter and his team have tried to mine the Internet for competitive technical intelligence (CTI). They tried to bring Research Profiling into Web resources

mining to discover competitive intelligence in commercial area through Google Soap Search API [3]. But it adopts statistics of search results rather than deep analysis of text. Therefore, we construct an integrated Research Profiling framework and a suit of technology methods from novel view, so as to depict a research field from multi-dimensions on the basis of web resources.

This demo reports one of the major outcomes of a project named Science Monitoring and Evaluation based on Scientific Web Resources (SMESWR), which is funded by National Key Technology R&D Program in the 11th Five Year Plan of China. The main goal of this project is to form a comprehensive methodology on automatically extracting intelligence from web resources especially for scientific research analysis. Developing technologies to detect scientific research activities, to monitor the progress of one research field, to track the evolution of one research topic or a research community, and to profile the key research unit is the heart of the project.

2. SMESWR Infrastructure

This framework consists of five major components: (1) *Web resource collecting*. We collect important institutional websites, news websites and newsgroups of related research fields, RSS from related website, OAI repository of one institution, personal homepages and blogs of one researcher from the Internet. (2) *Semantic knowledge extraction*. We named these extracted research terms and research objects as knowledge components, and model them based on the Research Ontology¹ which was proposed refer to the SWRC, we improved SWRC and re-organized some class in order to apply to our extraction task better. About the approach we choose to extract knowledge components is an integrated one which includes such as lexical-pattern approach and statistical approach. In fact, we choose some matured open source software such as GATE, Stanford Parser and KEA to provide basic NLP support, then we use some machine learning technologies to improve precise and we expand some relation rule set based on Hearst pattern for relation extraction. (3) *Knowledge repository construction*. Knowledge repository is composed of a series of extracted knowledge components with timestamps, for example, an instance of structure “class, research object, harvest time” is “research project, Science Monitoring and Evaluation based on scientific web resources, 2009-01-01”; Based on this computable data structure, knowledge repository is more clear and effective for future analysis. (4) *Semantic mining*. By

¹ <http://124.16.154.12/HotPortal/ResearchOntology.owl>

using a set of co-occurrence analysis and semantic mining methods, we perform semantic mining based on the data stored in knowledge repository, try to form panoramic perspective of a specific research field, perform burst detection, hot topic detection, timeline tracking and relation mining to discovery knowledge hidden behind the web resources. (5) *Research profiling*. In this process, we perform visualization analysis to profile the targeted research field. We detect the scientific research activities in one research field, figure out the key components in the area, depict the relation between those components, monitor the progress of one research field, track the evolution of one research topic or a research community.

Based on the Research Ontology mentioned above, we construct a knowledge repository containing all the classes, relationships and their instances. This knowledge repository is built upon RDF triples and then stored in MS SQL Server. It provides information retrieval, inference, and statistics interfaces to enable intelligent semantic queries and reasoning.

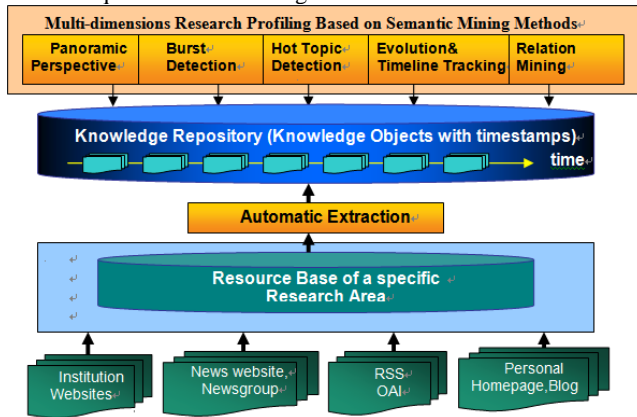


Figure 1. SMESWR Infrastructure

3. Implementation and Evaluation

At present, we choose "Artificial Intelligence" domain as our test domain. We have obtained 444 web site seed and 166 RSS harvest seeds. After the filter steps, we collected 89705 unique webpages and 18591 news articles. Based on these data, deep data mining has been conducted: through analyzing and scoring changes of each research object along the time axis, we identified the most important knowledge components in AI domain which provides the "big picture" of this field; tracking timeline, we use curve figures to illustrate the historical development of a certain knowledge object, and predict its future development trends; by analyzing the relationships among extracted knowledge components, we interlink the relation of different knowledge components by co-occurrence and relation extraction; and via clustering the top ranked research terms (e.g., top 2000 terms) along the timelines, the hot topics in AI have been visualized (Figure 2), more results displayed on the web portal*.

Comparing with related work, Arnetminer [4] mainly crawled academic personal websites and provides integrated overview of one researcher, but it does not conduct hot topic detection. CiteSeer [5] and GoogleScholar [6] are two of the largest collection of academic articles and provide basic citation data and

analysis, but it cannot monitor one specific organization or research field. Furthermore, all of them do not crawl news articles and take them into consideration for research evaluation.

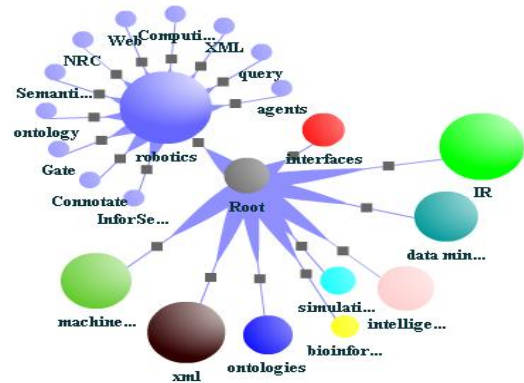


Figure 2. Hot Topic Detection

4. CONCLUSION

This paper describes a novel infrastructure to generate semantic-powered research profiling for research fields, organizations and individuals. It crawls related websites and news feeds, and on the basis of Research Ontology, they were modeled into RDF triples to facilitate semantic queries and semantic mining on burst detection, hot topic detection, dynamics of research, and relation mining based on the large-scale database. In the future, with the continuous accumulation of data, we hope to cluster the research terms periodically and tracking timeline of topic, and find topic changes through comparing every clustering result. Another important and hard-working task is to try to refine the extraction result to improve the performance. We also need to test the scalability and efficiency of our approaches and link our data with other Linked Open Data sets.

5. REFERENCES

- [1] Bollen, J., Rodriguez, M. A., & Van De Sompel, H. (2007). MESUR: Usage-based metrics of scholarly impact. In Proceedings of the 7th ACM/IEEE-CS joint conference on Digital libraries (Vancouver, BC, Canada, June 18-23, 2007), 474.
- [2] Carr, L., Bechhofer, S., Goble, C., & Hall, W. (2001). Conceptual linking: Ontology-based open hypermedia. In Proceedings of the 10th International WWW Conference (Hong Kong, China, May 01 - 05, 2001), 334-342.
- [3] Alan L. Porter, David J. Schoeneck, et al. 2007. Mining the Internet for Competitive Technical Intelligence. Competitive Intelligence Magazine 10 (May 5), 25-28
- [4] Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, & Zhong Su (2008). ArnetMiner: Extraction and Mining of Academic Social Networks. In Proceedings of the Fourteenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD'2008) (LAS VEGAS, US, Aug. 24-27) 990-998.
- [5] CiteSeer: <http://citeseer.ist.psu.edu/>
- [6] GoogleScholar: <http://scholar.google.com>

* <http://124.16.154.12/>