

Augmenting Lightweight Domain Ontologies with Social Evidence Sources

Albert Weichselbraun* · Gerhard Wohlgenannt \diamond · Arno Scharl \diamond

* Vienna University of Economics and Business
Department of Information Systems and Operations
Augasse 2-6, 1090 Vienna

`albert.weichselbraun@wu.ac.at`

\diamond MODUL University Vienna
Department of New Media Technology
Am Kahlenberg 1, 1190 Vienna

`gerhard.wohlgenannt@modul.ac.at`
`arno.scharl@modul.ac.at`

August 31, 2010

Agenda

Background and Motivation

Extracting Evidences from Social Sources

- Method

- Example Data

Evidence Integration

- System Diagram

- Spreading Activation

Evaluation

- Setting

- Informal Evaluation

- Formal Evaluation

Background and Motivation

- ▶ starting point: ontology learning framework (lightweight ontologies [Hendler, 2009, Alani et al., 2008])
- ▶ based on a seed ontology and domain documents
 - ▶ extract relevant terms
 - ▶ integrate them into the ontology
- ▶ benefits of integrating social sources
 - ▶ potential of providing background knowledge
 - ▶ contain the latest terminology [Angeletou et al., 2007] (evolve at much a higher pace as domain documents)

Extracting Evidences from Social Sources

- ▶ based on the seed terms \rightarrow transformation function (t) \rightarrow source specific (e.g. monograms for Delicious)
- ▶ disambiguation: WordNet
- ▶ social evidence sources:
 - ▶ easy Web Retrieval Toolkit (www.semanticlab.net/eWRT)
 - ▶ TagInfoService
 - ▶ implemented for Delicious (social bookmarking), flickr (photo/video hosting), technorati (blogs) and twitter (micro blogging)
- ▶ suggested tags \rightarrow relation weights based on the Dice coefficient

$$s_d(T_s, T_c) = \frac{2 \cdot n_{T_{sc}}}{n_{T_s} + n_{T_c}} \quad (1)$$

Example Triple Store Entries

seed ontology concept concept (C_s)	evidence source (e)	candidate concept (C_c)
climate change	oe:coOccurs	greenhouse gases
climate change	oe:twitterTag	environment
climate change	oe:deliciousTag	fuel

Table: Example evidence entries in the triple store.

Example Results

corpus-based	social	
	delicious	flickr
targets	animalcare	architecture
building	architects	art
coal	atmosphere	auckland
levels	award	beach
climate change policy	britney	bicycle
pact	carbonfootprint	brian
reduce greenhouse gas
pollution		
firm		
carbon dioxide emissions	technorati	twitter
ets	agile	aces
its carbon	apple	afghan
	architecture	afghanistan
	art	africa
	automotive	al_gore

Table: Terms from corpus-based and social evidence sources.

System Diagram

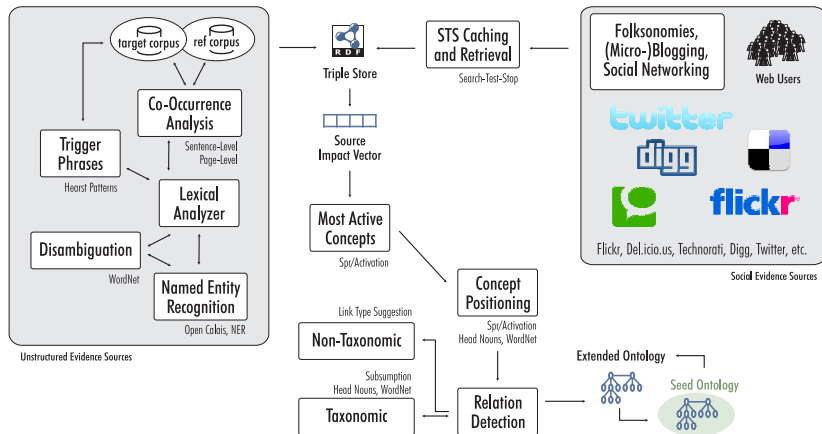


Figure: Ontology Extension Architecture System Diagram.

Spreading Activation

Goal: select the most promising candidate terms

Result from the previous process:

- ▶ evidence vector $\vec{r} \rightarrow$ contains evidence sources e :

$$\vec{r}(C_S, C_C) = \begin{pmatrix} r_{e_1}(C_S, C_C) \\ \dots \\ r_{e_n}(C_S, C_C) \end{pmatrix} \quad (2)$$

Transforming Evidences to Spreading Activation Weights

- ▶ Heuristic per-evidence-source translation rules s_e transform these relations using the source impact vector $\vec{S} = (s_{e_1}, s_{e_2}, \dots, s_{e_n})^T$ into a numerical weight

$$w(C_S, C_C) = |\vec{S}(\vec{r}(C_S, C_C))| \quad (3)$$

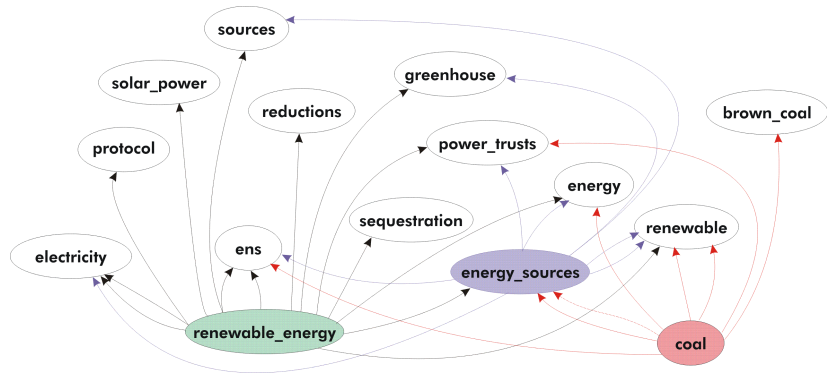
Evidence Integration - Example

$$\vec{r}(\text{cc}, \text{fuel}) = \begin{pmatrix} (oe : \text{coOccurs}, \text{sign} = 3.2) \\ (oe : \text{deliciousTag}, \text{dice} = 1.59) \\ (oe : \text{triggerPhrase}) \end{pmatrix}$$

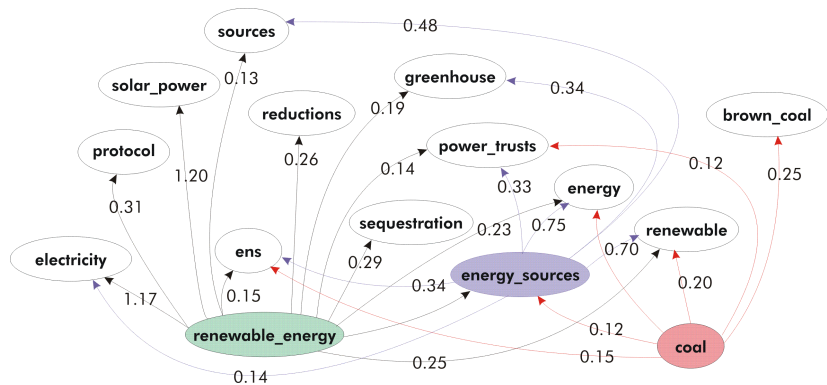
$$\vec{S} = \begin{pmatrix} 0.1 + 0.5 \cdot \text{sign} \\ 0.2 \cdot \text{dice} \\ 0.3 \end{pmatrix}$$

→ weight $w(\text{cc}, \text{fuel}) = 2.318$.

Evidence Integration - Spreading Activation Network



Evidence Integration - Spreading Activation Network



Evaluation - Setting

- ▶ seed “ontology”:
 - ▶ *fossil fuels* $\xrightarrow{\text{relatedTo}}$ *climate change* and
 - ▶ *fossil fuels* $\xrightarrow{\text{relatedTo}}$ *greenhouse gas(es)*
- ▶ domain corpora
 - ▶ 156 news media sites from the Newslink.org, Kidon.com and ABYZNewsLinks.com directories → 200,000 documents per week
 - ▶ six monthly corpora (April 2009 - August 2009)
 - ▶ domain detection based on regular expressions → climate change corpus containing 1250 documents / month
- ▶ social sources
 - ▶ Delicious, flickr, technorati, twitter
- ▶ two iterations → 24 new terms

Evaluation - Terms Removed and Added

terms removed	terms added
carbon dioxide emissions	agw
climate change policy	biomass
developing nations	cprs
kyoto protocol	cars
scientific assessments	epa
sulfur dioxide	ethanol
tom magliozzi	greenhouse-gas

Table: Selection of terms removed and added based on evidence from social sources.

Evaluation - Method

- ▶ pointwise mutual information (PMI)
 - how well are terms associated to each other
- ▶ four domain experts
 - relevance of the given relation
 - (0 .. irrelevant, 1 slightly relevant, 2 ...very relevant)

Evaluation - Pointwise Mutual Information

- ▶ Web metric (Yahoo! counts): seed tag counts (n_{T_s}), candidate tag (n_{T_c}) counts, common counts ($n_{T_{sc}}$)

$$n_z = n_{T_{sc}} + n_{T_s} + n_{T_c} \quad (4)$$

$$f(i) = \frac{n_i}{n_z} e^{-\frac{n_i}{n_z}} \quad (5)$$

$$PMI(T_s, T_c) = f(n_{T_{sc}}) / f(n_{T_s}) \cdot f(n_{T_c}) \quad (6)$$

Results

avg. PMI	corpus-based	corpus-based & social
April 2009	0.694 (16)	0.833 (17)
May 2009	0.753 (15)	0.921 (10)
June 2009	0.569 (16)	0.544 (15)
July 2009	0.625 (8)	0.862 (8)
August 2009	0.493 (5)	0.874 (9)
Sum	0.503 (60)	0.646 (59)




expert eval.	corpus-based	corpus-based & social
April 2009	0.875 (16)	1.353 (17)
May 2009	0.883 (15)	1.550 (10)
June 2009	1.000 (16)	1.283 (15)
July 2009	1.469 (8)	1.563 (8)
August 2009	1.150 (5)	1.167 (9)
Sum	1.013 (60)	1.369 (59)

Table: Impact of social evidence sources on ontology learning.

Outlook and Conclusions

- ▶ including social sources provides significant improvements to the ontology extension process (99.9% for a Welch two sample t-test and for the Wilcoxon rank sum test)
- ▶ drawbacks and potential pitfalls:
 - ▶ many social sources yield only unigrams
 - ▶ balancing corpus-based and social sources
- ▶ Future work:
 - ▶ support for n-grams
 - ▶ optimize source impact vectors based on user feedback
 - ▶ optimize access to remote resources (optimal stopping)

Thank you for your attention!

-  Alani, H., Hall, W., O'Hara, K., Shadbolt, N., Chandler, P., and Szomszor, M. (2008).
Building a pragmatic semantic web.
IEEE Intelligent Systems, 23(3):61–68.
-  Angeletou, S., Sabou, M., Specia, L., and Motta, E. (2007).
Bridging the gap between folksonomies and the semantic web:
An experience report.
In Workshop: Bridging the Gap between Semantic Web and Web, volume 2.
-  Hendler, J. (2009).
Web 3.0 emerging.
Computer, 42(1):111–113.