

Distributional Similarity of Multi-Word Expressions

Laura Ingram and James R. Curran

School of Information Technologies

The University of Sydney

NSW 2006, Australia

{ling6188, james}@it.usyd.edu.au

Abstract

Most existing systems for automatically extracting lexical-semantic resources neglect multi-word expressions (MWEs), even though approximately 30% of gold-standard thesauri entries are MWEs.

We present a distributional similarity system that identifies synonyms for MWEs. We extend Grefenstette's SEXTANT shallow parser to first identify bigram MWEs using collocation statistics from the Google WEB1T corpus. We extract contexts from WEB1T to increase coverage on the sparser bigrams.

1 Introduction

Lexical-semantic resources, such as WordNet (Fellbaum, 1998), are used in many applications in Natural Language Processing (NLP). Unfortunately, they are expensive and time-consuming to produce and are prone to bias and limited coverage. Automatically extracting these resources is crucial to overcoming the knowledge bottleneck in NLP.

Existing distributional approaches to semantic similarity focus on unigrams, with very little work on extracting synonyms for multi-word expressions (MWEs). In this work, we extend an existing system to support MWEs by identifying bigram MWEs using collocation statistics (Manning and Schütze, 1999). These are calculated using n-gram counts from the Google WEB1T corpus (Brants and Franz, 2006).

We evaluate against several gold-standard thesauri and observe a slight decrease in overall performance when the bigram MWEs were included. This is unsurprising since the larger vocabulary and sparser contextual information for bigrams makes

the task significantly harder. We also experimented with contexts extracted from WEB1T in an attempt to overcome the data sparseness problem. Inspection of the results for individual headwords revealed many cases where the synonyms returned were significantly better when bigram data was included.

2 Background

Distributional similarity relies on the *distributional hypothesis* that similar terms appear in similar contexts (Harris, 1954). Here we extend the SEXTANT parser (Grefenstette, 1994) to include multi-word *terms* and syntactic *contexts*.

Curran (2004) experiments with different parsers for extracting contextual information, including SEXTANT, MINIPAR (Lin, 1994), RASP (Briscoe and Carroll, 2002), and CASS (Abney, 1996). Lin (1998) used MINIPAR and Weeds (2003) used RASP for distributional similarity calculations. MINIPAR is the only parser to identify a range of MWEs. Weeds (2003) and Curran (2004) evaluate many measures for calculating distributional similarity. We follow (Curran, 2004) in using the weighted Jaccard measure with truncated *t*-test relation weighting for our experiments.

3 Detecting MWEs

The initial step in creating a thesaurus for MWEs is to identify potential MWE headwords using collocation statistics. We used various statistical tests, e.g. the *t*-test and the log-likelihood test (Manning and Schütze, 1999), calculated over the Google WEB1T unigram and bigram counts. These counts, calculated over 1 trillion words of web text, gave the most reliable counts. However, highly ranked terms, e.g.

Contact Us and Site Map, demonstrate bias towards website-related terminology. This list of known bigrams is used to detect bigrams within the BNC using the Viterbi algorithm.

4 Context Extraction

Grefenstette’s (2004) (SEXTANT) parser was extended to extract contextual information for the list of known known bigrams extracted above. Adding the knowledge of bigrams does not result in a substantial increase in the number of relations which implies that there is very little contextual information available about the bigram data. This has a significant impact on the difficulty of the task.

Experiments were also conducted whereby the contextual information was extracted from the WEB1T 4- and 5-gram data for a list of known bigrams from the gold-standard thesauri. This data lacks the syntactic information provided by SEXTANT but the counts are estimated over 10,000 times as much data. This should reduce the sparseness problem.

5 Synonym Extraction

Following Curran (2004), the proposed synonyms are compared directly against multiple gold-standard thesauri. We extend this evaluation to include multi-word headwords and synonyms. We randomly selected 300 unigram and 300 bigram headwords from the MAQUARIE (Bernard, 1990), MOBY (Ward, 1996), and ROGET’S (1911) thesauri, and WORDNET (Fellbaum, 1998).

We calculated the number of direct matches (DIRECT) and the inverse rank (INVR), the sum of the reciprocal ranks of matches. The results for the BNC-based experiments are summarised in Table 1.

Both INVR and DIRECT demonstrate the trend that performance decreases when MWES are included. However, performance did increase significantly for some terms when MWES were added. For example, tool improved from 0.270 to 0.627 INVR. The results for rate, shown in Table 2 also improved.

The next set of experiments extracted synonyms for 300 bigram headwords drawn from the MAQUARIE thesaurus. The best results for bigram headwords was achieved when unigram and bigram data was extracted from WEB1T and the VPC resource

DIRECT COUNT			
	UNIGRAM	BIGRAM	VPCS
TOTAL	6,731	6,623	6,286
AVERAGE	22.4	22.1	21.00
INVR			
	UNIGRAM	BIGRAM	VPCS
TOTAL	512.680	493.676	484.873
AVERAGE	1.709	1.646	1.616

Table 1: Comparison of DIRECT COUNT and INVR

UNIGRAMS	BIGRAMS	VPCS
level	level	level
price	price	price
cost	amount	cost
income	cost	amount
growth	speed	average

Table 2: Sample synonyms for rate

ATOMIC BOMB	DINING TABLE
atom bomb	coffee table
nuclear bomb	dining room
Atomic bomb	cocktail table
nuclear explosion	dining chair
Mushroom cloud	bedroom furniture

Table 3: Sample bigram synonyms

(Baldwin and Villavicencio, 2002) was included. Table 3 shows the top 5 synonyms (as ranked by the Jaccard measure) for atomic bomb and dining table.

6 Conclusion

We have integrated the identification of simple multi-word expressions (MWES) with a state-of-the-art distributional similarity system. We evaluated extracted synonyms for both unigram and bigram headwords against a gold-standard consisting of the union of multiple thesauri.

The main difficulties are the sparcity of distributional evidence for MWES and their low coverage in the gold standard. These preliminary experiments show the potential of distributional similarity for extracting lexical-semantic resources for both unigrams and MWES.

References

- Steven Abney. 1996. Partial parsing via finite-state cascades. *Journal of Natural Language Engineering*, 2(4):337–344.
- Timothy Baldwin and Aline Villavicencio. 2002. Extracting the unextractable: A case study on verb-particles. In *Proceedings of the Sixth Conference on Computational Natural Language Learning (CoNLL 2002)*, pages 98–104, Taipei, Taiwan.
- John R.L. Bernard, editor. 1990. *The Macquarie Encyclopedic Thesaurus*. The Macquarie Library, Sydney, Australia.
- Thorsten Brants and Alex Franz. 2006. Web1t 5-gram corpus version 1.1. Technical report, Google Research.
- Ted Briscoe and John Carroll. 2002. robust accurate statistical annotation of general text. In *Proceedings of the Third International Conference on Language Resources and Evaluation*, pages 1499–1504, Las Palmas de Gran Canaria, 29-31 May.
- James R. Curran. 2004. *From Distributional to Semantic Similarity*. Ph.D. thesis, University of Edinburgh.
- Chrisiane Fellbaum, editor. 1998. *Wordnet: An Electronic Lexical Database*. MIT Press.
- Gregory Grefenstette. 1994. *Explorations in Automatic Thesaurus Discovery*. Kluwer Academic Publishers, USA.
- Zellig Harris. 1954. Distributional structure. *Word*, 10(2/3):146–162.
- Dekang Lin. 1994. Principar - an efficient, broad-coverage, principle-based parser. In *Proceedings of COLING-94*, pages 482–488, Kyoto, Japan.
- Dekang Lin. 1998. Automatic retrieval and clustering of similar words. In *Proceedings of the 17th international conference on Computational linguistics*, volume 2, pages 768–774, Montreal, Quebec, Canada. Association for Computational Linguistics.
- Christopher D. Manning and Hinrich Schütze. 1999. *Foundations of Statistical Language Processing*. The MIT Press, London, England.
- Peter Mark Roget. 1911. *Thesaurus of English Words and Phrases*. Longmans, Green and Company, London, UK.
- Grady Ward. 1996. Moby thesaurus. <http://etext.icewire.com/moby/>.
- Julie E. Weeds. 2003. *Measures and Applications of Lexical Distributional Similarity*. Ph.D. thesis, University of Sussex.