

Integration of automatically-acquired multiword expressions in a hybrid machine translation system

Will Roberts and Markus Egg

Department of English and American Studies, Humboldt-Universität zu Berlin



Multiword expression (MWE) acquisition



- Unrestricted identification of senses by collecting lexical co-occurrence statistics on all words in Wikipedia
- Limited pre-processing of the text prior to sense identification:
 - Extract plain text from the Wikipedia dumps, segment text into sentences
 - tokenize and strip out text using regular expressions, removal of punctuation
 - No further processing (stemming, lemmatization, case normalization, removal of numbers or symbols)
 - Unstemmed text may be useful for capturing the morphological and syntactic features of some phenomena, such as e.g. will the team but not spill the beans
- Rank sense candidates using the log-likelihood association measure:
 - Collect word frequency information using the sense language modeling toolkit
 - Count n-grams with n up to 3 (i.e., we treat senses as bigrams and trigrams)

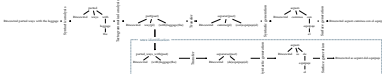
Compositionality ranking

- Take the top 10% from each association-measure-ranked list of senses and re-rank these candidates in order of increasing compositionality
- Based on Saha et al. (2007), this makes use of word embeddings, constructed using word2vec:
 - Build a vector representation for every word in the vocabulary as well as the every sense, using the extracted Wikipedia text
 - Gravely steep word-sense-embeddings of all occurrences of senses
 - Replace each of these with a single sense-word-sense index
 - Problem: gravely normalizing cannot handle senses which overlap
 - Solution: split senses into buckets with no overlaps
 - Each bucket produces a word-embedding space
 - Compare compositionality scores, and merge buckets back together
 - Compositionality score: cosine similarity of sense vector with the co-referent words (arithmetic mean)
 - Do not compare arbitrarily, with 'top sense' (the 50 most frequent senses in the vocabulary)

0.005 a forest for -0.005 -
 0.012 red tape -0.096 off -
 0.131 steps short of 0.215 a song -

Integration into the TechMT machine translation system (English-Spanish)

- TechMT (Zabokrievsky et al., 2008) is a hybrid machine translation system built on a pipeline model; statistical analysis phases (e.g., parsing, transfer) are interleaved with rule-based components
- The system analyses source text up to a high (hologrammatical) level of abstraction: a dependency graph containing only autonomous words.



sense identification is performed by string matching; successfully identified senses are collapsed into a single hologrammatical node.

Results

- QTLing test corpus contains 3K sentences, ca. 21K words of text from the TED domain
- sense scores for translation models trained on Europarl and in-domain (1.2M sentences, 1.2M words) test:

Europarl Sentences	Europarl Sentences			Test
	Language	Training	System	
baseline 23.24	28.60			
$\theta = 0.1$	24.21	28.60**	24.45**	4
$\theta = 0.2$	24.14	28.61**	24.41**	7
$\theta = 0.3$	24.11	28.61**	24.38**	8
$\theta = 0.4$	24.08	28.61**	24.36**	8
$\theta = 0.5$	24.06	28.61**	24.35**	8
*** - - - - -				
		28.60	24.45	4
		28.61	24.41	7
		28.61	24.38	8
		28.61	24.36	8
		28.61	24.35	8

Discussion

- Sense-only analysis of automatically acquired senses improves translation quality for this language pair (5% of 8% points)
- The improvement is only seen for the models built with the in-domain text:
 - An indication that our approach is sensitive to the domain of the training data
 - Evaluation paragraphs sensitive to the compositionality of the senses
 - The greatest improvements over the baseline are seen with small values of θ
 - Including more compositional senses ($\theta > 0.3$) eventually reduce word sense scores below the baseline
 - Composite models representing compositional senses likely cannot be adequately translated by single senses
- Methodology introduced here is:
 - Automatic and rule-free, allowing construction of linguistic resources with a minimum of human effort, requires no external lexical resources or language-specific tools
 - Language-independent
 - Domain-independent

References

- Saha, Indira, Paul Cook, and Timothy Baldwin. A word embedding approach to predicting the compositionality of multiword expressions. In *Proceedings of the 2011 Conference of the North American Chapter of the Association for Computational Linguistics*, pages 977-981, 2011.
- Zabokrievsky, Jan Prátek, and René Pösch. TechMT: Highly modular MT system with backgrounders used as transfer layer. In *Proceedings of the Third Workshop on Statistical Machine Translation*, pages 89-90. Association for Computational Linguistics, 2008.

Workshop on Statistical Machine Translation, 3-8 April 2010, Stroud, Massachusetts

