

WORD ASSISTANT

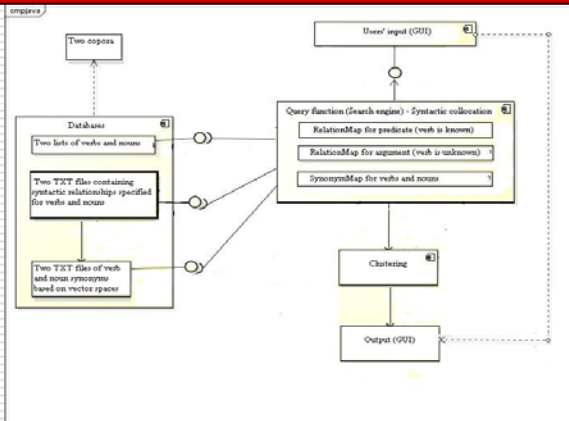
Zheng Ye, Philipp Busch, Pei Lu

Department of Computational Linguistics, University of Heidelberg, Germany

Introduction and Project Architecture

Problem: Due to limited vocabulary, non-native writers are often unsure about specific expressions.
Goal: The purpose of this software project is build a writing support system as one of possible solutions to the problem.
Architecture: The system consists of three modules, three databases, the query function and the GUI interface.

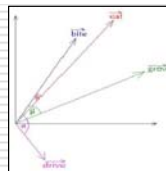
Process



Approaches

To achieve our goal, three special approaches are employed.

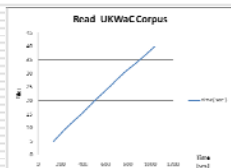
1. Three databases are extracted from the UK Web Archiving Consortium (UKWAC) and the British National Corpus (BNC). The UKWAC is a wide-ranged corpus covering diverse domains in English, while the BNC is manually annotated which avoids finding the wrong words.
2. The third database of synonyms is built with the help of two dependency vectors which are generated by DEPENDENCYVECTORS 2.5 (Pado 2010).



$$\text{sim}_{\cos}(\vec{x}, \vec{y}) = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

Evaluation

As it takes much time to read all files from the UKWAC and BNC, the software was tested with several mini examples. Herein is one of those examples.



Read resource

relationfre_noun.txt	dog love obj 7
relationfre_verb.txt	cat love obj 9
verbsynonymy.txt	dog love subj 5
nounsynonymy.txt	i love subj 10
verb_bnc.txt (same to the file in the resource)	you love obj 12
noun_bnc.txt (same to the file in the resource)	i want subj 9

Example	result
i # you	love want
Input: verb	
i love #	cat you dog
Input: noun	
# love dog	love you obj 11
Input: noun	love i subj 9
love ?you	love dog subj 10
Input: noun	love cat obj 16
i ?love you	want fun obj 12
Input: verb	want book obj 13
i love ?dog	want i subj 10
Input: noun	

3. A search engine is set up on the basis of two hash maps, i.e. RelationMap and SynonymMap, and uses frequency of occurrence of word combination to filter words.
4. The GUI is built as an interface to provide the user a direct access to Word Assistant with a dialog-box.

Conclusion

This software has the potential to be expanded despite its limitation on verbs and nouns. The experiment with verbs and nouns serves as a basis for future research. The obvious disadvantage is that it takes much time to read the files from the corpora due to their vast size.

References

A Ferraresi, E Zanchetta, M Baroni, S Bernardini (2008). *Introducing and evaluating ukWac, a very large web-derived corpus of English*. <http://clic.cimec.unitn.it/marco/publications/lrec2008/lrec08-ukwac.pdf> (Consulted on 15 January 2011).

S Pado, U Pado and K Erk (2010). A Flexible, Corpus-driven Model of Regular and Inverse Selectional Preferences. *Computational Linguistics* 36: 4.

S Pado (2010). *Manual for DependencyVectors 2.5*. http://www.nlpado.de/~sebastian/dv/dv_manual_2.5.pdf (Consulted on 20 December 2010).

philippbusch@gmx.net; pamelaz_yz@yahoo.de; lupei.hd@gmail.com