# Introducing a structure into a set of similar concepts

**Agnieszka Mykowiecka**[1,2] **and Małgorzata Marciniak**[1]

[1] Institute of Computer Science, Polish Academy of Sciences
Jana Kazimierza 5, 01-248 Warszawa,Poland
{agn,mm}@ipipan.waw.pl
[2] Polish-Japanese Academy of Information Technology
Koszykowa 86, 02-008 Warszawa, Poland

## Abstract

In the paper, we examine the idea of supporting domain ontology creation by an automatic clustering of selected terms identified using a terminology extraction method. We discuss the problem of introducing a structure into a set of similar concepts. We extract terminology from economic articles in Polish Wikipedia, then we select several sets of similar concepts present in the top 5,500 extracted terms. We describe two methods for automatic clustering of such groups of phrases on the basis of their distributional properties, i.e. the quantitative characteristics of the contexts of their occurrences in texts and test them on two sets of data.

## 1. Introduction

Ontologies consist of concepts organized in hierarchies thus building a domain ontology usually includes two steps: selecting a set of concepts that should be represented, and introducing relations which are held between them. To acquire a set of concept names for a domain we can extract a list of terms from a domain corpus. They can be manually organized into an ontology by experts, but the task can also be done with automatic support.

The paper addresses the problem of automatic identification of the multi-aspectual division of a subset of phrases which look like being co-hyponims, i.e. being sub-concepts of the one hiperonim. For example, the following concepts: *podatek liniowy* 'flat tax', *podatek progresywny* 'progressive tax', *podatek VAT* 'VAT tax' and *podatek od wynagrodzeń* 'payroll tax' are the sub-concepts of the concept *podatek* 'tax'. The first two refer to the way tax is counted, while the last two refer to the source of the amount being taxed. In this case, two different subcategorization criteria should be recognized and two sub-groups of the tax concept should be created, see Figure 1.
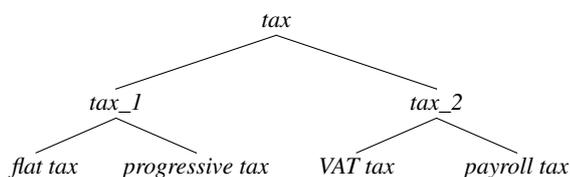


Figure 1: Hierarchy of taxes

In the paper, we describe how we select sets of sub-concepts and propose two methods for automatic clustering of these groups of phrases on the basis of the contexts of their occurrences in texts. The clustering methods are tested on two sets of data. The first one is small and of good quality. The second one consists of sentences from the Internet containing selected phrases. It is large, but it reflects the rather poor quality of the Internet texts.

The paper starts with a short description of the terminology extraction from the plWikiEcono corpus and the way groups of terms/concepts used in experiments are selected. Then, we describe how we collect and clean the large amount of data from the Internet. Finally, we describe the experiments with terms clustering, the results and evaluation.

## 2. Term extraction

Our approach to automatic terminology identification consists (like many others, see (Pazienza et al., 2005)) of two steps. The first one identifies candidates for terms on the basis of linguistic knowledge. Linguistic analysis of data starts with tokenization, morphology analysis and disambiguation. Then, nominal phrases are extracted using a cascade of shallow grammars. The grammar identifies the following constructions:

- single nouns or their equivalents;
- nouns followed or preceded by adjectival phrases;
- nouns followed by another noun in the genitive;
- combinations of the last two structures;
- noun phrases modified by prepositional phrases.

Candidates are ranked according to statistical information that indicates their importance in the analyzed texts. For ranking purposes, we use a slightly modified version of the C/NC method described in (Frantzi et al., 2000). In this method, all phrases are assigned a numerical value which is computed on the basis of the number of their occurrences within the text, the context in which they occur, and their length. Our modifications rely on one word terms being taken into account, and differentiating phrase contexts only on the basis of the neighboring words. The latter modification results in a slight change in ordering of terms (Marciniak and Mykowiecka, 2104). Terminology extraction was carried out on economic articles from Polish Wikipedia. From 5,500 top-ranked terms which occurred at least 3 times we select sets of sub-concepts that take part in the clustering experiments.

## 3. Selecting sets of sub-concepts

A set of sub-concepts can be selected from phrases which have the head element (the noun) related to the concept. In Polish, such phrases consist of:

- a noun and a modifying adjective such as *podatek$_{noun}$ dochodowy$_{adj}$* 'income tax';

- a noun and another noun in the genitive such as $prawo_{noun}$ $pracy_{noun,gen}$ 'labor law';
- a noun with a preposition modifier such as $podatek_{noun}$ $od_{prep}$ $wynagrodzeń_{noun}$ 'payroll tax'.

The above phrases are sub-concepts of *podatek* 'tax'. In the paper, we focus on sub-concepts consisting of a noun and its adjective modifier (the first item above).

Adjective modification is one of the most typical constructions occurring inside terminological phrases, and it has attracted some attention in the community working on automatic ontology learning. In (Almuhareb and Poesio, 2004) and (Cimiano, 2006), the authors proposed using adjectives for attribute learning. One of the subsequent works is (Hartung and Frank, 2010) in which a semi-supervised machine-learning approach for the classification of adjectives into property-denoting (like in *old box*) vs. relation denoting adjectives (like in *environmental science*) and object-denoting (like *eloquent person*) is presented (BEO classification). In addition to this classification, we observe term-denoting adjectives as being mostly part of lexicalized phrases. Although, initially, they might have represented relations, usually they are not separated from a described noun, and form a term or a name of a subtype of a concept. An example of such a phrase in Polish is $fundusz_{noun}$ $inwestycyjny_{adj}$ 'investment fund' which is likely to be represented as one concept not two concepts *fund* and *investment*. An even more evident example is $analiza_{noun}$ $techniczna_{adj}$ 'technical analysis' which is one term and will not be decomposed into either two concepts or a concept with an attribute. As differentiating these types of adjectives from relation-denoting adjectives could be difficult (for example, $podatek_{noun}$ $dochodowy_{adj}$ 'income tax' can be either represented as one concept or two related concepts depending on the desired level of description) we treat them both as one group.

The idea, which is widely accepted by linguists, is a division of adjectives into modifying and classifying ones based on the difference in their role within a noun phrase. To cluster noun phrases containing adjectives, first we have to filter out these adjectives which describe the properties of concepts, not the subtypes. In Polish, adjectives can occur on both sides of the noun. Adjectives preceding nouns (*duży wzrost* 'big increase') usually define features of a concept represented by a following noun phrase, so they are mostly property-denoting, while adjectives placed after a noun (*bank centralny* 'central bank'), have a classification role, thus they mostly belong to the relation-denoting class. The above observations lead to the conclusion that the main source of the important domain concept names are phrases with adjectives located to the right of a noun.

Within the 5500 top terms we identify noun and adjective sequences which constitute valid Polish phrases. This process results in obtaining 1500 nouns modified by 600 adjectives. Classification tests are performed on the subsets of this list containing 90 nouns which are modified by at least 4 different adjectives (345 adjectives in total). Only phrases which occur at least 5 times are taken into account in the clustering experiment. To evaluate the clustering methods, we randomly selected 11 nouns that create groups of sub-concepts. They consist of 5–23 phrases.

# 4. Data description

## 4.1. Wikipedia and supplementary data

The data set used in the study consisted of three segments coming from: economic articles from Polish Wikipedia (2.2 mln. tokens), a publicly available manually corrected subcorpus of the National Corpus of Polish (Przepiórkowski et al., 2012) (1.3 mln.) and a subset of texts from the newspaper "Rzeczpospolita" (0.5 mln.). The data was morphologically annotated using the Pantera tagger (Acedański, 2010). The NKJP subcorpus was manually corrected. The first set was used for the terminology extraction task and served as the source for concept names while the two others were used only an additional data for context frequencies. Table 1 gives the amount of selected phrase occurrences in the data.

Table 1: Phrase occurrences

| Noun | Translation | plWikiEcono | All |
|---|---|---|---|
| *prawo* | 'law' | 313 | 503 |
| *rynek* | 'market' | 293 | 418 |
| *fundusz* | 'fund' | 325 | 390 |
| *dochód* | 'income' | 135 | 147 |
| *podatek* | 'tax' | 313 | 384 |
| *spółka* | 'company' | 489 | 531 |
| *ekonomia* | 'economy' | 92 | 93 |
| *pieniądz* | 'money' | 43 | 50 |
| *jednostka* | 'unit/entity' | 134 | 203 |
| *grupa* | 'group' | 81 | 108 |
| *handel* | 'trade' | 88 | 104 |

## 4.2. Internet data

The second corpus consists of sentences containing selected phrases (sub-concept) collected from the Polish Internet.[1] The data was cleaned and corrected. We removed sentences containing strings longer than 30 characters. We took into account sentences containing 4–100 words, as shorter sentences do not contain enough context and longer ones are usually parts of tables or itemization. We corrected character coding and words divided by hyphenation. Finally, we removed duplicated sentences. The statistics of the cleaned data are given in Table 2. The data[2] was morphologically annotated using the same Pantera tagger.

# 5. Term similarity

The goal of the work was to recognize coherent subgroups of terms on the basis of term similarity. Two different ways of defining the likeness of terms were tested.

## 5.1. Hybrid similarity measure

The first similarity counting schema was established in a way that resembles (Nenadić et al., 2004). We count the similarity of terms based on their contexts, co-occurrence in sentences, and the similarity of adjectives being parts of these phrases. We also use information from Polish Wordnet version 2 (Piasecki et al., 2009).

---

[1] We want to express our gratitude to Dariusz Czerski for collecting the data.

[2] The data is available from `http://zil.ipipan.waw.pl/EconomicExcerpts/`.

Table 2: Internet data size

| Noun | Translation | Nb. of tokens | Nb. of phr. |
|------|-------------|---------------|-------------|
| *prawo* | 'law' | 22,818K | 763,851 |
| *rynek* | 'market' | 37,748K | 1,061,168 |
| *fundusz* | 'fund' | 22,695K | 700,958 |
| *dochód* | 'income' | 5,602K | 140,782 |
| *podatek* | 'tax' | 30,105K | 895,141 |
| *spółka* | 'company' | 37,748K | 1,061,168 |
| *ekonomia* | 'economy' | 4,591K | 154,582 |
| *pieniądz* | 'money' | 514K | 18,110 |
| *jednostka* | 'unit/entity' | 5,772K | 163,905 |
| *grupa* | 'group' | 30,796K | 903,558 |
| *handel* | 'trade' | 5,170K | 158,357 |

### 5.1.1. Partial similarity coefficients

**Contextual similarity** Contextual similarity is counted on the basis of tokens that appear around the term, without crossing sentence boundaries. In the case of the nearest base form we neglect the end of sentences or paragraphs, punctuation marks and conjunctions. All data below is separately counted for left and right contexts of terms:

- the sequences of 1 to 3 base forms;
- the sequences of POS tags of 2 to 4 tokens;
- the base form of the nearest verb;
- the base form of the nearest noun type token, adjectival token and preposition.

Contextual term similarity is counted separately for all of the features enumerated above using the Jaccard coefficient based on frequencies. We count the proportion of the common occurrences to the sum of all occurrences.

**Co-occurrence in a sentence.** In the paper (Nenadić et al., 2004) the *syntactical similarity* measure for pairs of terms is counted on the basis of their co-occurrence in patterns like: *such as*, *e.g.*, *like*, *both. . . and*, *including*, etc. No words except those explicitly cited in patterns and terms are allowed in the phrases. If we consider terms that consist of only two words, it is difficult to preserve such conditions. We observed that a sufficient similarity indication is if two terms with the same noun appear in the same sentence, so in the following sentence two terms in angle brackets represent concepts of the same type: *Unikatową, jak na <spółkę osobową>, cechą <spółki partnerskiej> jest . . .* 'As in a <partnership>, a unique feature of a <limited partnership> is . . .'.

**Common nouns** Many adjectives specify subtypes for many concepts, e.g. *sprawozdanie finansowe* 'financial report', *instrument finansowy* 'financial instrument', *piramida finansowa* 'financial pyramid'. If adjectives modify the same nouns, e.g. *krótkoterminowy* 'short-term' and *długoterminowy* 'longterm' modify 10 and 11 nouns respectively in Wikipedia economic texts, and 8 of them are the same, it may lead to the conclusion that they both describe one subcategorization criterion (in this case period of time). To account for this observation, we established a similarity measure whose value is equal to the proportion of the number of common nouns modified by two adjectives (from the two analyzed terms) to the maximum number of nouns modified by these two adjectives.

**Left and right descriptive adjectives within one noun phrase.** In Polish phrases where two classifying adjectives occur on both sides of the noun, one usually describes the phrase more precisely, e.g. $sejmowa_{adj}$ $komisja_n$ $budżetowa_{adj}$ 'parliamentary budget committee', but sometimes the second adjective gives information about other aspects of classification, e.g. $giełdowa_{adj}$ $spółka_n$ $odzieżowa_{adj}$ 'clothing trading company'. In these cases, the adjectives describe two different classification aspects. An adjective which is more important in a particular context is placed after a noun, so the order of adjectives can be changed according to context changes. To account for this fact, we assign a non-zero similarity value between phrases containing adjectives which in some cases surround one noun, and later we can use negative weight while combining this feature with others.

**Coordination of right modifiers.** In our data we have two types of phrases in which two adjectives (or in a very few cases more than 2) are located to the right of a noun. The first one is coordination, e.g. $komunikacja_n$ $kolejowa_{adj}$ $i_{conj}$ $lotnicza_{adj}$ 'rail and air transport'. The second type are constructions written with a hyphen, e.g. *problematyka społeczno-ekonomiczna* 'socio-economic problems'. These adjectival constructions represent the idea of something belonging to both groups A and B at the same time or being partially A and partially B. In both cases two classifying adjectives usually describe the same classification aspects.

**Wordnet similarity** As an additional information source we used plWordnet. It contains a very small number of multi word units (86% of tested phrases are not described), but all nouns and nearly all adjectives from this list (302 for 345 adjectives, 87.5%) are present. To calculate the similarity of our selected phrases we used information on three Wordnet relations: synonymy, antonymy and hyponymy, taking place between adjectives. As we did not have sense labeling, we used information on all senses described in plWordnet. For all tested phrases, only 25 pairs of adjectives modifying the same noun have a non-zero Wordnet similarity coefficient. Practically the only source of data was the antonymy relation.

### 5.1.2. Final (Overall) term similarity

The final similarity of a pair of terms was calculated as a weighted sum of 25 normalized coefficients:

- left/right POS contexts of length 2/3/4 (c2l, c3l, c4l, c2r, c3r, c4r);
- first left/right verb, noun, adjective, preposition (l_v, l_n, l_a, l_p, r_v, r_n, r_a, r_p);
- lemma right and left contexts of the length 1, 2 and 3 (bl_1, bl_2, bl_3, br_1, br_2, br_3);
- common noun modification coefficient (c-n);
- adjective co-occurrence in coordination (coord-adj) and two side modification use (s-adj);
- similarity: one common measure for contextual similarity and Wordnet similarity (sim);
- sentence term co-occurrence (sent).

### 5.2. Vector similarity measure

The second way to count similarity of terms is based on the standard cosine similarity measure, (Cimiano, 2006). Each phrase is represented by a vector, whose elements describe the occurrence of a particular word or a sequence of words in sentential contexts of the phrase. The features which are taken into account are based on word forms or lemmas and include:

- words from the both side context (*sg-win*);
- unigrams, bigrams and trigrams from the right and left windows, including skip-grams, i.e. n-grams consisting of non directly adjacent elements, (*ngr-win*);
- selected bigrams from a larger window, (*addbi-win*).

## 6. Clustering

A clustering experiment was performed on 11 selected sets of phrases (11 modified nouns). We asked two experts (with an economic background) to manually group phrases within the sets. The final grouping is a compromise between the two annotators. Table 4 shows the number of groups in the final division and the initial number of groups recognized by both annotators, together with the result of the point-wise F-measure (B-cubed measure (Bagga and Baldwin, 1998)). Then, for each set of phrases, automatic clustering was done using MultiDendrograms (Fernández and Gómez, 2008) performing hierarchical clustering on the basis of similarities counts described in Section 5. From several options available in the program, Ward's clustering algorithm was selected (Ward, 1963).

The results of automatic and manual clustering were compared using the B-cubed measure, see Table 5. For the first method of similarity counting, the results of two schema of combining similarities are given for both data sets. The first one consists in using uniform weights for all coefficients; in the second one, the weighting schemata of the coefficients was manually adjusted. This model (called *Man*) is given in Table 3. The tuning of parameters was done manually on the basis of the subjective judgment of divisions obtained for data not involved in the evaluation.

Table 3: The *Man* model

| coeff. | value | coeff. | value | coeff. | value | coeff. | value |
|---|---|---|---|---|---|---|---|
| left/right POS | | | | | | | |
| c2l | 0.05 | c3l | 0.05 | c4l | 0.01 | | |
| c2r | 0.05 | c3r | 0.05 | c4r | 0.01 | | |
| first verb, noun, adjective, prep | | | | | | | |
| l_v | 0.3 | l_n | 0.2 | l_a | 0.2 | l_p | 0.2 |
| r_v | 0.3 | r_n | 0.2 | r_a | 0.2 | r_p | 0.2 |
| lemma | | | | | | | |
| bl_1 | 0.3 | bl_2 | 0.2 | bl_3 | 0.2 | | |
| br_1 | 0.3 | br_2 | 0.1 | br_3 | 0.1 | | |
| c-n | 0.2 | sim | 0.2 | s-adj | 0.2 | | |
| coord-adj | 0.1 | sent | 0.1 | | | | |

For the second method, we show the results obtained for *sg-win*=6 and 10, *ngr-win*=5 and *addbi-win*=10 (named as 6-5-10 and 10-5-10). We observed that enlarging the window size lowered the results, i.e. for most cases the window of 10 elements was worse than the window of size 6. The results for larger windows were even worse. For

vector coordinates, we used either *tf-idf* or positive normalized *pmi* (Bouma, 2009) weights in place of frequencies. The better results were obtained for *npmi*. As an additional feature, we took bigrams for which the *npmi* counted for this particular text set was high (above 0.5) from the window of *addbi-win* size. We tested models based directly on word forms and on lemmas and observed that the vector based on forms gave better results in most cases, so we only give the results for these models. The last result shown in Table 5 was obtained for the continuous models built by word2vec (Mikolov et al., 2013) with the default values of parameters.

The first observation from the results included in Table 5 is that no method was shown to be clearly dominant. The word2vec models were better than any of our models for three out of eleven phrase sets. For our methods, the best results were the same for 4 phrase sets. For 3 sets, better results were obtained by the first method, while for 4 sets, by the second one. In the first method, the manual weights assignment improved the results, not only on the initial data set but also on the data set which was not inspected while establishing those weights. The results were worse for only one set (*unit*). Using much larger data lowered the results more frequently than it improved them and this decrease was more evident while using the first method. For 6 sets of phrases, the best result was obtained on the small Wikipedia-extended data, for 3 sets the better results were obtained on the Internet data, while for 2 sets, the F-measure was the same for both types of data.

Table 6 shows the best result obtained for the *market* group. In this case the agreement with manual clustering is high. One of the inconsistencies lies in recognizing two additional groups among the members of the fifth cluster which can be interpreted as more specific correct division. The only evident error is recognizing *agricultural market* as a member of the same group as the *internal market*.

Table 4: Manual clustering

| Noun | Phr. | Groups | An1 | An2 | F |
|---|---|---|---|---|---|
| law | 23 | 6 | 4 | 3 | 66.0 |
| market | 22 | 6 | 7 | 7 | 92,7 |
| fund | 17 | 5 | 3 | 6 | 55,9 |
| income | 14 | 7 | 7 | 7 | 82,7 |
| tax | 14 | 5 | 7 | 6 | 58,9 |
| company | 14 | 4 | 6 | 7 | 43,7 |
| economy | 12 | 2 | 4 | 3 | 66,0 |
| money | 7 | 2 | 2 | 2 | 100 |
| unit/entity | 7 | 2 | 2 | 2 | 100 |
| group | 6 | 5 | 3 | 4 | 87,5 |
| trade | 5 | 2 | 2 | 2 | 100 |

## 7. Conclusions

The obtained results confirm the great difficulty of the task of deciding whether or not a pair of terms belong to one aspect of subcategorization. The results also show that, among our methods, there is the lack of a model that clearly gives better results for all sets of phrases. However, although the results cannot be directly used within a new ontology, they can give additional information for ontol-

Table 5: Clustering results, different weights schemes and similarity measures

| | Wiki-ext | | | | Web | | | | Wiki-ext | | | | Web | | | | Web | |
| | Uniform | | Man | | Uniform | | Man | | 6-5-10 | | 10-5-10 | | 6-5-10 | | 10-5-10 | | word2vec | |
| | gr | F | gr | F | gr | F | gr | F | gr | F | gr | F | gr | F | gr | F | gr | F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| law | 13 | 46.1 | 10 | 50.7 | 6 | 53.2 | 10 | 57.1 | 10 | 69.1 | 11 | **72.3** | 8 | 68.3 | 4 | 62.6 | 10 | 69.5 |
| market | 10 | 61.0 | 13 | 62.2 | 7 | 59.6 | 12 | 69.6 | 6 | 74.0 | 6 | 74.0 | 8 | **78.7** | 10 | 72.1 | 6 | 77.0 |
| fund | 10 | 57.5 | 6 | **62.2** | 4 | 46.3 | 4 | 60.8 | 9 | 56.7 | 4 | 57.5 | 11 | 53.1 | 3 | **62.2** | 3 | **67.0** |
| income | 12 | 70.7 | 12 | 70.7 | 3 | 41.6 | 13 | 65.0 | 7 | 74.4 | 7 | **76.7** | 5 | 74.7 | 5 | 74.7 | 8 | 69.6 |
| tax | 7 | 67.6 | 9 | **73.5** | 7 | 50.2 | 13 | 58.6 | 9 | 67.8 | 10 | 69.6 | 10 | 62.8 | 8 | 62.9 | 4 | 67.6 |
| company | 6 | **66.9** | 3 | 65.9 | 6 | 59.3 | 5 | 59.3 | 6 | 57.7 | 6 | 55.9 | 5 | 56.9 | 6 | 63.1 | 4 | 57.3 |
| economy | 2 | 48.9 | 2 | 54.1 | 4 | 46.5 | 2 | 54.1 | 2 | 67.7 | 2 | 66.7 | 2 | 65.7 | 2 | **68.8** | 2 | **73.7** |
| money | 4 | 59.6 | 2 | **80.7** | 2 | 68.5 | 2 | 68.5 | 2 | 65.7 | 2 | 65.7 | 2 | 65.7 | 3 | 57.1 | 2 | 55.2 |
| unit/entity | 3 | **88.0** | 2 | 78.6 | 3 | 68.6 | 4 | 74.4 | 3 | **88.0** | 3 | 83.3 | 2 | 68.5 | 2 | 68.5 | 3 | 83.3 |
| group | 3 | 71.8 | 5 | **83.3** | 2 | 69.6 | 5 | **83.3** | 5 | **83.3** | 5 | **83.3** | 5 | **83.3** | 5 | **83.3** | 3 | **84.6** |
| trade | 4 | 69.6 | 4 | 78.7 | 2 | **100.** | 2 | **100.** | 2 | **100.** | 2 | **100.** | 4 | 75.0 | 4 | 84.6 | 2 | **100.** |

Table 6: The results for *rynek* 'market'

| Manual clustering | Automatic |
|---|---|
| rynek docelowy 'target market' | A |
| rynek nowy 'new market' | A |
| rynek efektywny 'efficient market' | B |
| rynek rolny 'agricultural market' | C |
| rynek wewnętrzny 'internal market' | C |
| rynek krajowy 'domestic market' | D |
| rynek międzynarodowy 'international market' | D |
| rynek zagraniczny 'foreign market' | D |
| rynek światowy 'worldwide market' | D |
| rynek lokalny 'local market' | D |
| rynek finansowy 'finacial market' | E |
| rynek kapitałowy 'capital market' | E |
| rynek kredytowy 'credid market' | E |
| rynek pieniężny 'money market' | E |
| rynek ubezpieczeniowy 'insurance market' | E |
| rynek walutowy 'currency market' | E |
| rynek giełdowy 'exhange' | F |
| rynek równoległy 'parallel market' | F |
| rynek kasowy 'cash (spot) market' | G |
| rynek terminowy 'futures market' | G |
| rynek pierwotny 'primary market' | H |
| rynek wtórny 'secondary market' | H |

ogy creators as to whether or not two concepts should be treated as co-hyponims or rather belong to different sub-categorization dimensions.

The next interesting conclusion is that the large amount of Internet data does not give definitely better results than small encyclopedic texts in which the number of phrases occurrences is relatively low. To confirm this observation, it would be interesting to carry out more experiments concerning the other models and types of phrases.

## 8. References

Acedański, Sz., 2010. A morphosyntactic brill tagger for inflectional languages. In H. Loftsson, E. Rögnvaldsson, and S. Helgadóttir (eds.), *Advances in Natural Language Processing*, volume 6233 of *LNCS*. Springer.

Almuhareb, A. and M. Poesio, 2004. Attribute-based and value-based clustering. an evaluation. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*.

Bagga, A. and B. Baldwin, 1998. Algorithms for scoring coreference chains. In *LREC Workshop on Linguistics Coreference*.

Bouma, G., 2009. Normalized (pointwise) mutual information in collocation extraction. In *From Form to Meaning: Processing Texts Automatically, Proceedings of the Biennial GSCL Conference 2009*. Tübingen.

Cimiano, P., 2006. *Ontology Learning and Population from Text. Algorithms, Evaluation and Applications*. Springer.

Fernández, A. and S. Gómez, 2008. Solving non-uniqueness in agglomerative hierarchical clustering using multidendrograms. *Journal of Classification*, 25:43–65.

Frantzi, K., S. Ananiadou, and H. Mima, 2000. Automatic recognition of multi-word terms: the C-value/NC-value method. *Int. Journal on Digital Libraries*, 3:115–130.

Hartung, M. and A. Frank, 2010. A semi-supervised type-based classification of adjectives: Distinguishing properties and relations. In *Proc. of LREC 2010*. ELRA.

Marciniak, M. and A. Mykowiecka, 2104. Terminology extraction from medical texts in Polish. *Journal of Biomedical Semantics*, 5:24.

Mikolov, T., K. Chen, G. Corrado, and J. Dean, 2013. Efficient estimation of word representations in vector space. In *Proceedings of Workshop at ICLR*.

Nenadić, G., I. Spasić, and S. Ananiadou, 2004. Automatic discovery of term similarities using pattern mining. *International Journal of Terminology*, 10(1):55–80.

Pazienza, M., M. Pennacchiotti, and F. Zanzotto, 2005. Terminology Extraction: An Analysis of Linguistic and Statistical Approaches. In S. Sirmakessis (ed.), *Knowledge Mining Series: Studies in Fuzziness and Soft Computing*. Springer Verlag.

Piasecki, M., S. Szpakowicz, and B. Broda, 2009. *A Wordnet from the Ground Up*. Wrocław: Oficyna Wydawnicza Politechniki Wrocławskiej.

Przepiórkowski, A., M. Bańko, R. Górski, and B. Lewandowska Tomaszczyk (eds.), 2012. *Narodowy Korpus Języka Polskiego*. Wydawnictwo Naukowe PWN.

Ward, Joe H., 1963. Hierarchical grouping to optimize an objective function. *Journal of the American Statistical Association*, 58(301):236–244.