

Automatic Word Clustering in Studying Semantic Structure of Texts

Olga Mitrofanova

St. Petersburg State University, Faculty of Philology and Arts,
Department of Mathematical Linguistics,
Universitetskaya emb., 11
199034 St. Petersburg, Russia
{alkonost-om@yandex.ru}

Abstract. The purpose of the study is to prove that results of automatic word clustering (AWC) may contribute much in investigating semantic structure of texts and in evaluating plot complexity. Experiments were carried out for Russian texts, mainly stories and short novels. Data obtained in course of study allowed to formulate and verify several linguistic hypotheses.

Keywords: Automatic Word Clustering, Russian Corpora, Semantic Structure of Texts

1 Introduction

Formalization of text structure and quantitative evaluation of semantic relations between text units prove to be of considerable importance in various fields of natural language understanding: modelling plot structure, text summarization, evaluation of translation adequacy in parallel texts, automatic text indexing, classification of texts in corpora, etc. (for a detailed analysis cf. [1], [2]).

One of the procedures providing linguistic data on semantic structure of texts is automatic word clustering (AWC). It is assumed that AWC results help to reveal semantic structure of texts and to determine plot complexity. To prove this assumption, AWC procedure was carried out with the help of a specialized AWC toolkit based on word space model. Experimental procedure implied processing Russian texts, mainly stories and short novels. A set of key words describing major topics of the plot was assigned to each text, clusters of words with similar distributions were created for each key word. Data extracted from texts through AWC procedure admit thorough linguistic interpretation. Further comparison of cluster content and structure allowed to distinguish texts characterized by a plot including a dominating topic with a number of subtopics and texts characterized by a plot including a set of major (independent or correlating) topics.

© A. Gelbukh (Ed.)
Advances in Computational Linguistics.
Research in Computing Science 41, 2009, pp. 27-34

Received 07/12/08
Accepted 11/12/08
Final version 05/02/09

2 AWC Procedure

From a linguistic point of view, AWC is based on the possibility of detecting semantic similarity of words by comparing their syntagmatic properties (co-occurrence or distribution analysis); from a technical standpoint, AWC involves construction of vector-space models for processed texts; it means that the sets of contexts for each word are represented as distribution vectors in N -dimensional space [5], [7].

It is possible to evaluate semantic similarity of words by measuring distances between their distribution representations. Numerous metrics are used for the given purpose. The selection of metrics often depends on qualitative parameters of processed texts. In our case, cosine measure (*Cos*) was chosen as a basic metric. Results of measuring semantic distances are applied in clustering: words having similar distribution representations as a rule reveal similarity of meaning and should be included into the same cluster.

General approaches to clustering are exposed in hierarchical (agglomerative, divisive), partitioning (K -means, K -medoid, etc.), hybrid algorithms. Certain linguistic tasks require application of special clustering techniques, e.g. CBC [4], MajorClust [6], etc. The choice of a particular algorithm is determined by experimental conditions (corpora size, required speed of clustering, constraints for the size of resulting clusters, etc.). In our research preference was given to agglomerative clustering algorithm as it seems to be applicable in case of limited data and appropriate for processing texts of small / medium size.

Experiments were carried out with the help of AWC tool [3]. Python-based AWC software maintains procedures of text preprocessing and agglomerative clustering. Such parameters as names of input files (processed texts and key words describing the content of a text), context window size, weight assignment for context items, size of clusters, etc. are determined by users.

Text preprocessing is performed at the first stage. Context segmentation is carried out in accordance with a particular context window size. Automatic weight assignment may be done for lexical items taking into account their positions in contexts. Then, distribution representations of words are formed, co-occurrence matrix is built, semantic distances are calculated at the second step. These data are necessary for agglomerative clustering which is performed at the third step. An output file contains clusters of words with similar distributions in a text, such clusters being formed for each key word.

3 Linguistic Data

Experiments were carried out for over 20 Russian texts, mainly stories and short novels (cf. table 1). The texts differ in authorship (A. Belyaev, M. Bulgakov, N. Gogol, A. Grin, E. Zamyatin, A. Žitinsky, etc.), in size ($N = 8\,491 \dots 37\,217$ tokens), in lexical diversity (number of unique words $L = 3\,038 \dots 6\,144$ tokens, Somers coefficient $S = \ln \ln L / \ln \ln N = 0.920 \dots 0.936$). In some experiments both raw and morphologically tagged texts were subjected to analysis. Processing raw texts provides data on distribution of word forms (tokens), while processing morpholo-

gically tagged texts allows to reveal interrelations between words (lemmas) within texts. In particular cases original Russian texts and their translations were considered as well. The texts were extracted from M. Moškov digital library (<http://lib.ru/>). Frequency lists for each text were created, additional statistical information (frequencies of words from various parts of speech, average sentence length, amount of dialogues, etc.) was obtained with the help of FantLab linguistic processor (<http://www.fantlab.ru/>).

Table 1. Texts subjected to analysis.

Author, title	Size (tokens) (N)	Number of unique words (L)	Somers coefficient (S)
Gogol N. <i>Taras Bul'ba</i>	37 217	6 144	0.920
Žitinsky A. <i>Časy s variantami</i> (<i>A Clock with Variants</i>)	28 092	5 197	0.922
Belyaev A. <i>Poslednij čelovek iz Atlantidy</i> (<i>The Last Man of Atlantis</i>)	26 892	5 160	0.924
Bulgakov M. <i>Sobačje serdce</i> (<i>Dog's Heart</i>)	25 218	5 321	0.928
Bulgakov M. <i>Rokovyje jajca</i> (<i>The Fatal Eggs</i>)	21 199	5 084	0.933
Zamyatin E. <i>Na kuličkah</i> (<i>In Kulički</i>)	20 832	4 544	0.928
Grin A. <i>Alyje parusa</i> (<i>Crimson Sails</i>)	20 366	4 984	0.933
Grin A. <i>Priklučeniya Ginča</i> (<i>Ginč's Adventures</i>)	19 120	5 017	0.936
Belyaev A. <i>Večny hleb</i> (<i>Eternal Bread</i>)	17 103	3 640	0.924
Grin A. <i>Kolonija Lanfier</i> (<i>Lanfier Colony</i>)	15 532	3 943	0.932
Belyaev A. <i>Mertvaja golova</i> (<i>A Dead Head</i>)	14 820	3 519	0.928
Gogol N. <i>Povest' o tom, kak possorilis' Ivan Ivanovič s Ivanom Nikiforovičem</i> (<i>A Tale of How Ivan Ivanovič Quarrelled with Ivan Nikiforovič</i>)	14 052	3 071	0.923
Belyaev A. <i>Zolotaja gora</i> (<i>A Golgen Hill</i>)	12 505	3 008	0.927
Belyaev A. <i>Čelovek, kotoryj ne spit</i> (<i>A Sleepless Man</i>)	11 943	3 104	0.931
Gogol N. <i>Viy</i>	11 800	2 824	0.926
Bulgakov M. <i>Zapisky na manžetah</i> (<i>Notes on the Cuff</i>)	10 056	3 038	0.937
Belyaev A. <i>Ni žizn', ni sm'ert'</i> (<i>Neither Life nor Death</i>)	9 681	2 653	0.931
Bulgakov M. <i>Morphij</i> (<i>Morphia</i>)	8 491	2 493	0.934

4 Experimental Results

In course of experiments a set of five key words – frequent words describing major topics of the plot – was assigned to each text, e.g.:

Zamyatin E. *Na kuličkah* (*In Kulički*):

key words {*kapitan* (captain), *Tihmen'*, *Marus'a*, *Andrej*, *Šmit*};

Žitinsky A. *Časy s variantami* (*A Clock with Variants*):

key words {*žizn'* (life), *vrem'a* (time), *časny* (watch), *ded* (grandfather), *ja* (I)}.

Clusters of lexical items with similar distributions were created for each key word. The following parameters of clustering were chosen in the experiments: similarity measure – *Cos*, context window size – ± 5 , size of clusters – 10 items, no weight assignment. Previously it was found out that AWC performed with such parameters provides quite reliable data. Resulting clusters contain words or word forms associated with key words in a text and ordered according to *Cos* values. Distances between key words and their nearest neighbours in clusters (*D*) and difference between D_{\max} and D_{\min} in clusters (*Var*) were calculated for each text.

Table 2. Example (1): clusters of word forms extracted for key words in texts.

Text:	Bulgakov	M.	Morphij	(Morphia);
key words	<i>Polyakov, doktor (doctor),</i>	<i>otdelenije (department),</i>	<i>pis'mo (letter),</i>	<i>Marja;</i>
cluster elements are ordered in accordance with <i>Cos</i> values				
<i>Polyakov</i>	<i>doktor (doctor)</i>	<i>otdelenije (department)</i>	<i>pis'mo (letter)</i>	<i>Marja</i>
<i>pipiska (postscript)</i> 0.328	<i>krasu (beauty)</i> 0.390	<i>terapevtičeskoje</i>	<i>nelepoje (absurd)</i> 0.428	<i>Vlasjevna</i> 0.731
<i>krupnymi (large)</i> 0.293	<i>zastrelils'a</i>	<i>(therapeutic)</i> 0.589	<i>isteričeskoje (hysterical)</i> 0.349	<i>prolepetala (prattled)</i> 0.326
<i>bukvami (letters)</i> 0.293	<i>(shot himself)</i> 0.390	<i>doktoru (doctor)</i> 0.490	153 0.349	<i>divženije (movement)</i> 0.320
<i>smerti (death)</i> 0.289	<i>užas (horror)</i> 0.388	<i>hirurgičeskoje (surgery)</i> 0.431	<i>sarkoma (sarcoma)</i> 0.309	<i>šlepnula (slapped)</i> 0.320
<i>umer (died)</i> 0.285	<i>takoj (such)</i> 0.388	<i>Pavlu (Paul)</i> 0.423	<i>duše (soul)</i> 0.299	<i>bormotala (muttered)</i> 0.320
<i>krasu (beauty)</i> 0.254	<i>jehala (drove)</i> 0.379	<i>zaraznoje (infectious)</i> 0.409	<i>načalo (beginning)</i> 0.295	<i>brauning (Browning)</i> 0.287
<i>pomutneli (dimmed)</i> 0.219	<i>umer (died)</i> 0.333	<i>deckoje (infant)</i> 0.382	<i>roždalos' (was borning)</i> 0.284	<i>zadeta (touched)</i> 0.281
<i>mimoletmuju (fleeting)</i> 0.219	<i>drožala (trembled)</i> 0.323	<i>akušerskoje (obstetric)</i> 0.374	<i>ležalo (lay)</i> 0.259	<i>cepko (firmly)</i> 0.281
<i>slyšno (audible)</i> 0.215	<i>doroga (road)</i> 0.231	<i>mašina (car)</i> 0.340	<i>razdražat' (annoy)</i> 0.259	<i>boleznetno (painfully)</i> 0.281
	<i>lampoj (lamp)</i> 0.231	<i>bol'soj (big)</i> 0.272		

Table 3. Example (2): clusters of word forms extracted for key words in texts.

Text:	Belyaev	A.	Čelovek, kotoryj ne spit	(A	Sleepless	Man);
key	word	preparat	preparat	(medicine),	(medicine),	
cluster elements are ordered in accordance with <i>Cos</i> values						
preparat (medicine)						
	<i>himiki (chemists)</i> 0.259					
	<i>gotovyj (ready)</i> 0.259					
	<i>prodažu (sale)</i> 0.236					
	<i>uničtožavšij (destroying)</i> 0.233					
	<i>polučils'a (came out)</i> 0.227					
	<i>obnaružili (discovered)</i> 0.227					
	<i>polipeptidy (polypeptides)</i> 0.195					
	<i>vypuskalo (produced)</i> 0.169					
	<i>najdeny (found)</i> 0.163					

It seems that cluster elements often correspond to essential features of objects, persons or events denoted by key words and somehow emphasized in a text.

Relations between cluster elements can be characterized as syntagmatic and / or paradigmatic, e.g. synonymy & attributive relation: *terapevtičeskoje (therapeutic), hirurgičeskoje (surgery), zaraznoje (infectious), deckoje (infant), akušerskoje (obstetric) – otdelenije (department)*; meronymy: *otdelenije (department) – doktor (doctor)*; person – actions: *Marja – prolepetala (prattled), šlepnula (slapped), bormotala (muttered)*, phraseological units and compounds: *Marja – Vlasjevna* (first name & second name), etc. (cf. table 2).

Those relations can be properly described in terms of semantic roles and lexical functions, e.g. action *obnaružili (discovered)* – agent *himiki (chemists)*, result

preparat (medicine) – attribute *gotovij (ready)*, *uničtožavšij (destroying)*; action *prodažu (sale)* – theme *preparat (medicine)*, etc. (cf. table 3).

Thus, AWC allows to reveal and analyze not only standard but also occasional relations between lexical items which may be specific for a particular text or a set of texts of the same author or dealing with the same topic.

In some tests clustering was performed in two modes: with weight assignment and without weight assignment. In most cases clusters contain similar elements – word forms (tokens) in raw texts or words (lemmas) in tagged texts. At the same time those words or word forms within clusters may be ordered differently as regards their *Cos* values. So, clusters may be similar in content, but they may differ in structure (cf. table 4). It should be noted that in experiments with weight assignment *Cos* values for nearest neighbours of key words in clusters (*D*) seem to be lower than in experiments without weight assignment.

Table 4. Example: clusters obtained in experiments with / without weight assignment.

Text: Gogol N. <i>Vij</i>; key word <i>bursak (seminarist)</i>, cluster elements are ordered in accordance with <i>Cos</i> values	
Clustering without weight assignment	Clustering with weight assignment
<i>bursak (seminarist)</i>	<i>bursak (seminarist)</i>
<i>sodrognuls'a (shuddered) 0.479</i>	<i>sodrognuls'a (shuddered) 0.436</i>
<i>pozelenevšije (green) 0.442</i>	<i>pozelenevšije (green) 0.405</i>
<i>ostupivši (having stepped aside) 0.420</i>	<i>holod (cold) 0.371</i>
<i>vperil (stared) 0.379</i>	<i>izumlenija (amusement) 0.364</i>
<i>holod (cold) 0.359</i>	<i>mertvyje (dead) 0.338</i>
<i>čuvstvitel'no (perceptibly) 0.299</i>	<i>čuvstvitel'no (perceptibly) 0.305</i>
<i>izumlenija (amusement) 0.295</i>	<i>sv'atoj (saint) 0.222</i>
<i>žizni (life) 0.259</i>	<i>probežal (run) 0.205</i>
<i>sv'atoj (saint) 0.200</i>	<i>žizni (life) 0.176</i>

We also considered clustering results obtained in course of processing raw texts and morphologically tagged texts. Correspondence of word forms (tokens) and words (lemmas) in clusters created for raw and morphologically tagged texts (cf. table 5) proves the existence of stable intrinsic relations underlying text structure. These relations remain almost intact as the analysis moves from the level of word forms (tokens) to the level of words (lemmas). So, AWC procedure may furnish us with additional information on the integrity and continuity of the text as a complex of heterogeneous linguistic units.

AWC proves to be of much use in comparative analysis of original texts and translations, as it often allows to evaluate stylistic and semantic similarity of texts. Similarity of clusters formed for a word and its translation equivalent reveals correspondence between contexts of those words in the original and in translation, while differences of content and structure of such clusters imply syntactic / morphological / lexical differences of texts in question as well as inconsistency in the choice of translation equivalents for a particular word or for lexical items co-occurring with this word in contexts (cf. table 6).

As statistical parameters of texts may influence results of clustering, additional tests were required. We've studied the texts written by A. Belyaev which reveal

common semantic structure and are characterized by a branching plot with numerous and frequently changing topics. The given texts differ in size and in number of unique words. At the same time, distances between key words and their nearest neighbours in clusters don't vary much for those texts ($D \in [0.088 \dots 0.259]$). It turns out that such parameters as size and number of unique words play important but not decisive role in studying text structure by means of AWC.

Table 5. Example: clusters obtained in experiments with raw and tagged texts.

Text: Bestužev-Marlinsky A. <i>Strašnoje gadanje (A Scary Fortune-telling)</i> ;	
key	word
cluster elements are ordered in accordance with <i>Cos</i> values	
Clustering in a raw text (tokens)	Clustering in a tagged text (lemmas)
<i>neznakomec (stranger)</i>	<i>neznakomec (stranger)</i>
<i>stenky (wall)</i> 0.219	<i>drognut' (quaver)</i> 0.223
<i>podjezdu (entrance)</i> 0.219	<i>trost' (cane)</i> 0.198
<i>vysadiv (having put off)</i> 0.219	<i>vysadit' (pull off)</i> 0.197
<i>večor (evening)</i> 0.218	<i>zajti (overstep)</i> 0.196
<i>zahvatyvaja (seizing)</i> 0.216	<i>večor (evening)</i> 0.196
<i>rasseržen (angry)</i> 0.216	<i>kalitka (gate)</i> 0.195
<i>trost' (cane)</i> 0.193	<i>zahvatyvut' (seize)</i> 0.194
<i>gorst'ami (in handfuls)</i> 0.188	<i>rasserdit' (anger)</i> 0.192
<i>car'a (tzar)</i> 0.172	<i>ironičeskij (ironical)</i> 0.173
<i>ironičeskoju (ironical)</i> 0.165	<i>gorst' (in handfuls)</i> 0.169

Table 6. Example: comparison of clusters formed for test words in the original text and in translation.

Texts:	Grin	A.	Alyje	parusa	(Crimson	Sails);
test	words	<i>Sekret</i>		(<i>Secret</i>),	<i>galiot</i>	(<i>galliot</i>),
cluster elements are ordered in accordance with <i>Cos</i> values						
Russian text	English text	Russian text	English text	Russian text	English text	
<i>Sekret (Secret)</i>	<i>Secret</i>	<i>Galiot (Galliot)</i>	<i>galliot</i>			
<i>potr'asenija (shock)</i> 0.239	<i>intimations</i> 0.213	<i>trehmačtovyj (three-mastered)</i> 0.800	<i>masted</i> 0.843			
<i>vdohnovenogo (inspired)</i> 0.210	<i>hurries</i> 0.212	<i>dvesti (two hundred)</i> 0.700	<i>purchased</i> 0.556			
<i>dvesti (two hundred)</i> 0.178	<i>agitation</i> 0.202	<i>šest 'des' at (sixty)</i> 0.600	<i>sixty</i> 0.527			
<i>neuderžimymi (uncontrollable)</i> 0.178	<i>rounding</i> 0.201	<i>kuplennyj (purchased)</i> 0.600	<i>ton</i> 0.509			
<i>slezami (tears)</i> 0.178	<i>shock</i> 0.173	<i>tonn (ton)</i> 0.500	<i>brig</i> 0.316			
<i>mravits'a (likes)</i> 0.149	<i>cape</i> 0.173	<i>Grejem (Gray)</i> 0.329	<i>hundred</i> 0.271			
<i>kamenistoj (rocky)</i> 0.149	<i>uncontrollable</i> 0.144	<i>sobstvennikom (proprietor)</i> 0.3	<i>orion</i> 0.222			
<i>padajuš'im (falling)</i> 0.147	<i>masted</i> 0.117	<i>kapitanom (captain)</i> 0.291	<i>rugged</i> 0.211			
<i>golovokžitel'no (astounding)</i> 0.117	<i>galliot</i> 0.093	<i>mačty (masts)</i> 0.290	<i>Arthur</i> 0.189			

Thorough treatment of AWC results allowed us to distinguish three main types of texts with regard to their semantic structure (Types 1, 2, and 3).

Type 1 is represented by texts characterized by a plot including a dominating topic with a number of subtopics. For such texts distances between key words and their nearest neighbours in clusters (D) and difference between D_{\max} and D_{\min} in clusters (Var) are as follows: $D \geq 0.300$, $Var \geq 0.200$.

Type 2 is represented by texts characterized by a plot including a set of major (probably independent) topics. For such texts distances between key words and their nearest neighbours in clusters (D) and difference between D_{\max} and D_{\min} in clusters (Var) are as follows: $D < 0.300$, $Var < 0.200$.

Type 3 is represented by texts characterized by a plot including a set of major (probably correlating) topics. For such texts distances between key words and their nearest neighbours in clusters (D) and difference between D_{\max} and D_{\min} in clusters (Var) are as follows: $D \geq 0.300$, $Var < 0.200$.

Examples of texts representing Types 1, 2, and 3 are given in table 7.

Table 7. Texts representing Types 1, 2 and 3.

Type, author, title	D	Var
Type 1		
Gogol N. <i>Taras Bul'ba</i>	0.379	0.252
Grin A. <i>Priklučenija Ginča (Ginč's Adventures)</i>	0.406	0.231
Gogol N. <i>Povest' o tom, kak possorilis' Ivan Ivanovič s Ivanov Nikiforovičem (A Tale of How Ivan Ivanovič Quarrelled with Ivan Nikiforovič)</i>	0.453	0.357
Gogol N. <i>Viy</i>	0.547	0.424
Belyaev A. <i>Zolotaja gora (A Golden Hill)</i>	0.566	0.471
Bulgakov M. <i>Morphij (Morphia)</i>	0.731	0.403
Type 2		
Žitinsky A. <i>Časy s variantami (A Clock with Variants)</i>	0.149	0.048
Zamyatin E. <i>Na kuličkah (In Kulički)</i>	0.174	0.068
Grin A. <i>Alyje parusa (Crimson Sails)</i>	0.204	0.091
Bulgakov M. <i>Sobačje serdce (Dog's Heart)</i>	0.212	0.103
Belyaev A. <i>Ni žizn', ni sm'ert' (Neither Life nor Death)</i>	0.222	0.070
Belyaev A. <i>Poslednij čelovek iz Atlantidy (The Last Man of Atlantis)</i>	0.224	0.115
Grin A. <i>Kolonija Lanfier (Lanfier Colony)</i>	0.224	0.139
Belyaev A. <i>Mertvaja golova (A Dead Head)</i>	0.243	0.155
Belyaev A. <i>Čelovek, kotoryj ne spit (A Sleepless Man)</i>	0.259	0.144
Belyaev A. <i>Eternal bread (Večny Hleb)</i>	0.268	0.130
Bulgakov M. <i>Rokovyje jajca (The Fatal Eggs)</i>	0.279	0.182
Type 3		
Bulgakov M. <i>Zapisky na manžetah (Notes on the Cuff)</i>	0.359	0.091

Our observations on semantic structure of texts require more detailed consideration and further verification.

5 Conclusion

In course of our experiments performed for Russian stories and short novels we proved that AWC may be of great help in distinguishing three types of texts as regards their semantic structure. We managed to describe texts of different plot complexity: texts revealing a dominating topic and a set of subtopics, texts revealing a set of major (probably independent) topics, and texts revealing a set of major (probably correlating) topics. Linguistic analysis of cluster content and structure

allowed to study standard as well as occasional semantic relations between lexical items occurring in texts. Experiments on AWC performed for raw and morphologically tagged texts proved the existence of intrinsic relations underlying text structure, those relations being preserved at two levels of analysis: the level of word forms (tokens) and the level of words (lemmas). Comparison of AWC results obtained for the original texts and their translations proved to be relevant in the evaluation of stylistic and semantic similarity of texts.

Further research implies experiments carried out for texts of different size, genre and authorship, with expanded sets of key words, with changing parameters (context window size, cluster size, etc.).

Acknowledgements. The author expresses deep gratitude to CICLing–2009 Program and Organizing Committees for having accepted the paper, to reviewers for their kind attention and wise comments which helped to improve the paper, to Mikhail Alexandrov (Barcelona, Spain) and Elena Iagounova (St. Petersburg, Russia) for their deep interest in the given research and for their help and encouragement.

References

1. Bolshakov, I.A., Gelbukh, A.: Computational Linguistics: Models, Resources, Applications. IPN – UNAM – Fondo de Cultura Económica (2004)
2. Leontjeva, N.N.: Avtomatičeskoje Ponimanije Tekstov: Sistemy, Modeli, Resursy. Moscow (2006)
3. Mitrofanova, O., Mukhin, A., Panicheva, P., Savitsky, V.: Automatic Word Clustering in Russian Texts. In: Matoušek, V., Mautner, P. et al. (eds.): Text, Speech and Dialogue. Proceedings of the Tenth International Conference TSD–2007, Pilsen, Czech Republic, September 3–7, 2007. Lecture Notes in Artificial Intelligence, Vol. 4629. Springer-Verlag, Berlin Heidelberg New York (2007) 85–91
4. Pantel, P.: Clustering by Committee. Ph.D. Dissertation, Department of Computing Science, University of Alberta (2003) <http://www.isi.edu/~pantel/Content/publications.htm>
5. Sahlgren, M.: The Word-Space Model: Using Distributional Analysis to Represent Syntagmatic and Paradigmatic Relations between Words in High-dimensional Vector Spaces. Ph.D. dissertation, Department of Linguistics, Stockholm University (2006) <http://www.sics.se/~mange/TheWordSpaceModel.pdf>
6. Stein, B., Meyer zu Eissen, S.: Document Categorization with MajorClust. In: Proceedings of the 12th Workshop on Information Technology and Systems WITS–02. Barcelona, Spain (2002) 91–96
7. Widdows, D.: Geometry and Meaning. Center for the Study of Language and Information – Lecture Notes, Vol. 172. The University of Chicago Press (2004)