



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Information Processing and Management xxx (2006) xxx–xxx

**INFORMATION
PROCESSING
&
MANAGEMENT**www.elsevier.com/locate/infoproman

Text mining without document context

Eric SanJuan ^a, Fidelia Ibekwe-SanJuan ^{b,*}^a *LITA, Université de Metz & URI, INIST-CNRS, Ile du Saulcy, 57045 Metz Cedex 1, France*^b *URSIDOC-ENSSIB & Université de Lyon 3, 4, cours Albert Thomas, 69008 Lyon Cedex, France*

Received 16 March 2006; accepted 16 March 2006

8 Abstract

9 We consider a challenging clustering task: the clustering of multi-word terms without document co-occurrence infor-
10 mation in order to form coherent groups of topics. For this task, we developed a methodology taking as input multi-word
11 terms and lexico-syntactic relations between them. Our clustering algorithm, named CPCL is implemented in the Term-
12 Watch system. We compared CPCL to other existing clustering algorithms, namely hierarchical and partitioning (*k*-means,
13 *k*-medoids). This out-of-context clustering task led us to adapt multi-word term representation for statistical methods and
14 also to refine an existing cluster evaluation metric, the editing distance in order to evaluate the methods. Evaluation was
15 carried out on a list of multi-word terms from the genomic field which comes with a hand built taxonomy. Results showed
16 that while *k*-means and *k*-medoids obtained good scores on the editing distance, they were very sensitive to term length.
17 CPCL on the other hand obtained a better cluster homogeneity score and was less sensitive to term length. Also, CPCL
18 showed good adaptability for handling very large and sparse matrices.
19 © 2006 Published by Elsevier Ltd.

20 *Keywords:* Multi-word term clustering; Lexico-syntactic relations; Text mining; Informetrics; Cluster evaluation

22 1. Introduction

23 We developed a fast and efficient text mining system that builds clusters of noun phrases (multi-word terms)
24 without need of document co-occurrence information. This is useful for mapping out research topics at the
25 micro-level. Because we do not consider the within document co-occurrence, our approach can be conceived
26 as an *out-of-context clustering* except if we consider the *intraterm* context, i.e., words appearing in the same
27 terms can be said to share a similar context. Terms are clustered depending on the presence and number of
28 shared linguistic relations. For instance, a link will be established between the two terms *humoral immune*
29 *response* and *humoral Bhx immune response* since one is lexically included in the other. Likewise *clustering algo-*
30 *rithm* is linked to *computer algorithm* by a modifier substitution. This lexico-syntactic approach is suitable for
31 clustering multi-word text units which rarely re-occur as is in the texts. Such multi-word terms (MWTs) often

* Corresponding author.

E-mail addresses: eric.sanjuan@univ-metz.fr (E. SanJuan), ibekwe@univ-lyon3.fr (F. Ibekwe-SanJuan).

32 result in very large and sparse matrices or graphs¹ that are difficult to handle by the existing approaches to
33 clustering which rely on high frequency information. The resulting system, called TermWatch (Ibekwe-San
34 Juan, 1998a; SanJuan, Dowdall, Ibekwe-SanJuan, & Rinaldi, 2005) can be applied to several tasks like domain
35 topic mapping, text mining, query refinement or question–answering (Q–A).

36 Some attempts have been made to cluster document contents in the bibliometrics, scientometrics and infor-
37 metrics fields. Some authors have considered the clustering of keywords, classification codes or subject head-
38 ings assigned to documents by indexers (Braam et al., 1991; Callon, Courtial, & Laville, 1991; Zitt &
39 Bassecoulard, 1994). Although these information units depict the thematic contents of documents, they are
40 external to the documents themselves and do not allow for a fine-grained analysis of the current topics
41 addressed in the full texts. In studies where the document contents were considered, only lone words were
42 extracted through statistical analysis. The majority of clustering methods used in the information retrieval
43 field (Cutting, Karger, Pedersen, & Tukey, 1992; Eisen, Spellman, Brown, & Botstein, 1998; Karypis, Han,
44 & Kumar, 1994) are also based on the vector-space representation model of documents (bag-of-words
45 approach). To reduce the dimensions of the vector space, words with a discriminating power are selected based
46 on term weighting indices like the *Inverse Document Frequency* (IDF), *Mutual Information* (MI) or the cosine
47 measure. This also results in the drastic elimination of more than half of the initial data from the analysis. Our
48 text mining approach treats highly frequent and low frequent terms equally. This is important for applications
49 like science and technology watch where the focus is on novel information often characterised by low fre-
50 quency units (weak signals). Price and Thelwall (2005) have demonstrated the usefulness of low frequency
51 words for scientific Web intelligence (SWI). They showed that removing low frequency words reduced cluster
52 coherence and separation, i.e., clusters were less dissimilar.

53 Glenisson, Glänzel, Janssens, and Moor (2005) proposed combining full text analysis with bibliometric
54 analysis in order to cluster the research themes of 85 scientific papers. Text contents were represented as vec-
55 tors of lone words. Stemming was performed on the words and bigrams were detected, i.e. sequences of two
56 adjacent words that occurred frequently. It is a well-known fact that stemming brutally removes the semantics
57 of derived or inflected words. For instance, “stationary, station, stationed” are all reduced to *station*. Also,
58 bigrams may not always correspond to valid domain terms. The authors weighted the bigrams using the Dun-
59 ning likelihood ratio test (Dunning, 1993). This led to selecting the 500 topmost bigrams for analysis and dis-
60 carding the rest. One of the interesting findings of this study is that clustering items from full texts rather than
61 keywords or terms from the reference section leads to a more fine-grained and accurate mapping of research
62 topics. This finding is in line with our text mining approach.

63 Polanco, Grivel, and Royauté (1995) developed the Stanalyst informetrics platform. Stanalyst comprises a
64 linguistic component which identifies variants of MWTs used to augment their occurrences. The MWTs are
65 then clustered based on document co-occurrence information. To the best of our knowledge, no informetric
66 method has considered clustering phrases based on linguistic relations. The TermWatch approach is based on
67 the hypothesis that clustering multiword terms (MWTs) through lexico-syntactic and semantic relations can
68 yield meaningful clusters for various applications. In view of this, we developed a methodology that can han-
69 dle very large and sparse matrices in real time. For instance, in the current experiment, the input list of terms is
70 31,398, none which is eliminated prior to the matrix reduction phase.

71 The clustering algorithm implemented in TermWatch is named CPCL (Classification by Preferential Clus-
72 tered Link). This algorithm was first published in (Ibekwe-SanJuan, 1998a) but owing to its fundamental dif-
73 ferences with existing approaches, setting up an adequate comparison framework with other methods has been
74 a bottleneck issue. In this paper, we focus on the evaluation with other clustering algorithms (variants of par-
75 titioning and hierarchical algorithms). Evaluation is carried out on a test corpus (the GENIA project) which
76 comes with an answer key (gold standard). This will ensure that the results being presented are grounded in the
77 real world.

78 The rest of the paper is organised as follows: Section 2 gives details of the test corpus; Section 3 describes
79 our text mining methodology; Section 4 presents the evaluation method; Section 5 describes the experimental

¹ In the experiments run up to date, we have been able to handle graphs of 80,000 terms in real time applications for online data analysis and query refinement.

80 setup; Section 6 discusses the results of the evaluation with other clustering methods; Section 7 draws remarks
81 and conclusions.

82 2. Test corpus

83 In order to carry out an evaluation, we chose a dataset with an existing *ideal partition* (gold standard). The
84 GENIA project² consists of 2000 abstracts downloaded from the MEDLINE database using the search key-
85 words: *Human*, *Blood Cells*, and *Transcription Factors*. Biologists manually annotated the valid domain terms
86 in these texts, yielding 31,398 terms. This ensures in our experiment that competing methods start from the
87 same input. The GENIA project also furnished a hand-built ontology, i.e. a hierarchy of these domain terms
88 arranged into semantic categories. There are 36 such categories at the leaf nodes. Each term in the GENIA
89 corpus was assigned a semantic category at the leaf node of the ontology. We shall refer to the leaf node cat-
90 egories as *classes* henceforth. Of course, the GENIA ontology's hierarchy, the number of classes and the
91 semantic category of each term were hidden from the clustering methods. It should be noted that since the
92 GENIA ontology is a result of a human semantic and pragmatic analysis, we do not expect automatic clus-
93 tering methods to reproduce it exactly without prior and adequate semantic knowledge. The goal of the eval-
94 uation is to determine the method whose output requires the least effort to reproduce the classes at the leaf
95 nodes of the ontology. Also, it is worth noting that although the authors of this project use the term *ontology*
96 to qualify this hierarchy, it is more of a small taxonomy. Indeed, the GENIA *ontology* is still embryonic
97 because of its small size (36 classes, 31,398 terms). The classes are of varying sizes. The largest class, called
98 *other name* has 10,505 terms followed by the *protein molecule* class with 3899 terms and the *dna domain or*
99 *region* class with 3677 terms. The 12 smallest classes (*rna domain or region inorganic*, *rna substructure*, *nucle-*
100 *otide*, *atom*, *dna substructure*, *mono cell*, *rna nla*, *protein nla*, *carbohydrate*, *dna nla*, *protein substructure*) each
101 has less than 100 terms. It is quite revealing that the largest class is a miscellaneous class. This suggests that
102 this class can be further refined. Also some relations normally found in a full-fledged ontology are absent (syn-
103 onymy in particular). This tends to suggest that this hierarchy is a weaker semantic structure than an ontology
104 and can thus constitute an adequate clustering task. For these reasons, we prefer to refer to it as the GENIA
105 *taxonomy* henceforth.

106 Table 1 gives some examples of terms in the GENIA corpus.

107 Fig. 1 shows the fast decreasing distribution of terms in the 35 classes. We omitted the largest class, called
108 *other name* which concentrated 33% of the terms because it was difficult to fit in. A few number of classes (*pro-*
109 *tein molecule*, *dna domain or region*, *protein family or group*, *cell line*, *cell type*) concentrated the rest of the
110 terms (almost 75%). The bars show the proportion of terms according to their length. As a consequence of
111 this fast decreasing model, a clustering method optimised for one of the prominent classes can obtain good
112 scores without correctly classifying terms in the majority of the smaller classes. Another feature that can be
113 observed in Fig. 1 is that the distribution of one word terms is not correlated with the general distribution
114 of terms. Meanwhile, we will see in Section 6 that most of the clustering methods perform better on long terms
115 and thus on classes like “*protein family or group*” and “*dna domain or region*” that contain few one word
116 terms. In an OTC task, the intrinsic properties of MWTs (like term length) obviously play an important role
117 since they are the only available context.

118 3. Overview of our text mining methodology

119 Our methodology consists of three major components: MWT extraction; relation identifier and clustering
120 module. An integrated visualisation package³ can be used if topic mapping is the target. In this experiment,
121 this aspect will not be explored as evaluation will focus on cluster quality and not on their layout. However,
122 interested readers can find an application of research topic mapping in (Ibekwe-SanJuan & SanJuan, 2004).

² <http://www-tsujii.is.s.u-tokyo.ac.jp/GENIA/>.

³ The aiSee visualization package (<http://www.aisee.com>) has been integrated to the system.

Table 1
Examples of terms in GENIA corpus

GENIA category	Terms
Amino acid monomer	Amide-containing amino acid Asparagines <i>n</i> -Acetylcysteine
Atom	Cytosolic calcium feca2+
Body part	Organ Peripheral lymphoid organ Tumor-draining lymph node
Cell component	1389 sites/cell b6d2f1 mouse uterine cytosol Cytoplasmic protein extract il-13-treated human peripheral monocyte nuclear extract
Cell line	Anergized <i>t</i> cell Adherence-isolated monocyte Xenopus hepatocyte
Other name	Anatomic tumor size Apoptosis Follicular lymphoma

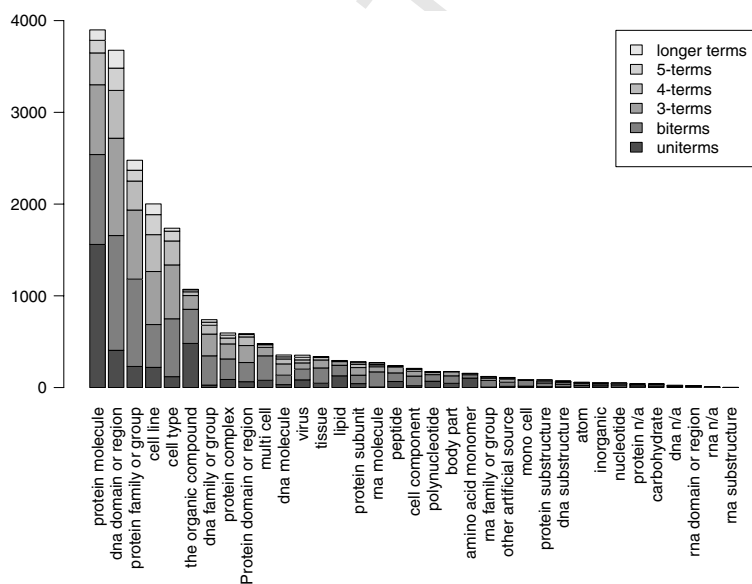


Fig. 1. Distribution of terms in GENIA categories.

123 3.1. Term extraction module

124 Note that in the current experiment, our term extraction module was not used as the terms were already
 125 manually annotated in the corpus. We however describe summarily its principle. TermWatch performs term
 126 extraction based on shallow natural language processing (NLP) techniques. Extraction is implemented via the
 127 NLP package developed by the University of Edinburgh. LTPOS is a probabilistic part-of-speech tagger based
 128 on Hidden Markov Models. It uses the Penn Treebank tag set which ensures the portability of the tagged texts
 129 with many other systems. LTCHUNK identifies simplex noun phrases (NPs), i.e., NPs without prepositional

130 attachments. In order to extract more complex terms, we wrote contextual rules to identify complex termino-
131 logical NPs. An example is provided in [Appendix A](#). About 10 such contextual rules were sufficient to take
132 care of the different syntactic structures in which nominal terms appear in English. Given that some domain
133 concepts can appear as long sequences like in *parental granulocyte-macrophage colony-stimulating factor (GM-
134 CSF)-dependent cell line*, it is obvious that such MWTs are not likely to re-occur frequently in the corpus.
135 Hence, the difficulty of clustering them with methods based on co-occurrence criteria.

136 3.2. Relation identifier

137 Different linguistic operations can occur within NPs. These operations either modify the structure or
138 the length of an existing term. They have come to be known as *variations* and have been well studied
139 in the computational terminology field ([Ibekwe-SanJuan, 1998b](#); [Jacquemin, 2001](#)). Variations occur at differ-
140 ent linguistic levels: morphological (gender and spelling variants), lexical (substitution of one word by another
141 in an existing term), syntactic (expansion or structural transformation of a term), semantic (synonyms, gen-
142 eric/specific relations). Our relation identifier tries to acquire all these types of variations among the input
143 terms.

144 3.2.1. Morphological variants

145 These refer to number (*tumor cell nuclei/tumor cell nucleus*) and gender variations in a term and also to
146 spelling variants. They enable us to recognise different appearances of the same term. For instance, *IL-9-
147 induced cell proliferation* will be recognised as a spelling variant of *IL 9-induced cell proliferation*. Spelling vari-
148 ants are identified using cues such as special characters while gender and number variants are identified using
149 WordNet ([Fellbaum, 1998](#)).

150 3.2.2. Lexical variants

151 We call substitution variants operations involving the change of only one word in a term, either in the mod-
152 ifier position (*coronary heart disease* ↔ *coronary lung disease*) or in the head position (*mutant motif* ↔ *mutant
153 strain*). The head is the noun focus in an English NP, i.e., the subject while the modifier plays the role of a
154 qualifier (an adjective). The head word is usually the last noun in a compound phrase (*strain in mutant strain*)
155 or the last noun before a preposition in a prepositional structure (*retrieval in retrieval of information*).

156 3.2.3. Syntactic variants

157 These refer to the addition of one or more words to an existing term as in *information retrieval and efficient
158 retrieval of information*. We call these operations *expansions*. Expansions that affect the modifier words are fur-
159 ther broken down into left-expansion and insertion. Alternatively, expansions can affect the head word. In this
160 case, we talk of *right expansion*.

161 Morphological variants (spelling) and permutation variants are recognised first since they refer to the same
162 term. Then these variants are used to recognise the more complex variants. For instance, *B cell development
163* haven been recognised as a spelling variant of *B-cell development*, this enables the identification of other types
164 of variants (syntactic and lexical) containing the two spelling variants. Variations are assigned a role during
165 clustering depending on their interpretation. This will be further detailed in Section 5.

166 3.2.4. Semantic variants

167 It is an accepted fact that syntactic relations suggest semantic ones (left expansions and insertion can engen-
168 der *generic-specific* links, some substitution variants can reflect see also relations). However, these semantic
169 relations are not explicit. Moreover, the types of relations considered so far all require one stringent condition:
170 that the related terms share some common words. This leaves out terms which can be semantically-linked but
171 without sharing common words, i.e. synonyms. In order to acquire explicit semantic links, we need an external
172 semantic resource. For this purpose, we chose WordNet ([Fellbaum, 1998](#)), a large coverage semantic database
173 which organises English words into synsets. A synset is a particular sense of a given word. Since WordNet
174 organises only words and not multi-word terms, we had to devise rules in order to map *word-word* semantic
175 relations into “*MWT–MWT*” relations in our corpus. One way to achieve this is to replace words by their

176 synsets and then apply the same variation relations to sequence of synsets. However, like all external
 177 resources, WordNet has some limitations. First is its incompleteness vis-à-vis specialised domain terminology.
 178 Second, being a general purpose semantic database, WordNet establishes links which can be incorrect in a
 179 specialised domain.

180 We thus restricted the use of WordNet to filtering out lexical substitutions, and consequently to pairs of
 181 terms that share at least one word in order to reduce the number of wrong semantic links. Only a very few
 182 number of relations were found. The following rule was applied to lexical substitutions in order to identify
 183 the semantic ones using WordNet hierarchy: given two terms related by a lexical substitution, check if the
 184 two words substituted are linked by an ascending or descending path in the hierarchy. Observe that, by def-
 185 inition of lexical substitutions, this rule only applies to words that are in the same grammatical position (head
 186 or modifier).

187 In this way, we acquired the following synonymy relations:

T cell growth ~ *T cell maturation*

189 *antenatal steroid treatment* ~ *prenatal steroid treatment*

190 Only 365 WordNet modifier substitutions and 208 WordNet head substitutions were found whereas lexico-
 191 syntactic variants were much more abundant (see Table 2).

192 Table 2 gives the number of variants identified for each type among the GENIA terms. As a term can be
 193 related to many others, the number of relations is always higher than the number of terms.

194 Details of the variation identification rules are given in Appendix B.

195 3.3. Clustering module

196 The TermWatch system implements a graph-based approach of the hierarchical clustering called CPCL
 197 (Classification by Preferential Clustered Link) originally introduced by Ibekwe-SanJuan (1998b). The main
 198 features of this approach are

- 199 (1) the intuitiveness of its results for human users since any pair of terms clustered together are related by a
 200 relative short path of real linguistic relations,
- 201 (2) an ultrametric model that ensures the existence of a unique and robust solution,
- 202 (3) its linear time complexity on the number of variations that allows interactive data analysis since cluster-
 203 ing can be processed in real time.

204
 205 We show here that this algorithm can be applied to other types of inputs. For that, we need to cast the
 206 description of the algorithm in the more general context of data analysis.

207 Let S be a sparse similarity data matrix defined on a set Ω of objects. This matrix can be represented advan-
 208 tageously by a valuated graph $G = (\Omega, E, s)$ where E is the set of edges made of all unordered pairs $\{i, j\}$ of
 209 objects such that $S_{ij} > 0$ and s is the valuation of edges defined for all $(i, j) \in \Omega^2$ by $s(i, j) = S_{ij}$. In the case
 210 of sparse data, the size of E is much smaller than $|\Omega|^2$.

Table 2
 Statistics on variation relations per type

Variation relation	Terms	Relations
Spelling variants	1560	2442
Left Right-expansions (exp_2)	294	441
Right-expansions (exp_r)	2329	3501
Left-expansions (exp_l)	2818	4260
Insertions (ins)	526	798
Modifier-substitutions (sub_mod.)	4291	3773
Head-substitutions (sub_head)	781	1082
WordNet-synonyms (sub_wn)	365	208

211 Let Val_S be the set of values in S . If $|Val_S| \ll |S|$ then, the usual hierarchical algorithms will produce small
 212 dendrograms since they will have at most $|Val_S|$ levels. Thus, they will induce a very reduced number of inter-
 213 mediary balanced partitions in the gap between the trivial discrete partition and the family of connected com-
 214 ponents of G . A way to correct this drawback of hierarchical clustering without losing its intuitiveness and
 215 computer tractability is not to consider edge values in an absolute way but in the context of adjacent edges.
 216 Thus, two objects related by an edge e will be clustered at a given iteration, only if the value of e is greater than
 217 any other value in its neighborhood. This means that i, j will be clustered at the first iteration only if S_{ij} is
 218 greater than the maximum in the line S_i and in the column S_j . It has been shown in [Berry, Kaba, Nadif,
 219 SanJuan, and Sigayret \(2004\)](#) that this variant of hierarchical clustering preserves its main ultrametric
 220 properties.

221 This solution is specially well adapted when the observed similarities between objects are generated by pair-
 222 wise observations. In the case of out-of-context clustering (OTC), given three terms u, v, t such that v shares at
 223 least one word with u and t (possibly not the same), we will consider a local criteria to decide if v is closest to u
 224 or to t .

225 In this approach, the clustering phase can be easily implemented using the following straightforward proce-
 226 dure which we call *SLME* (Select Local Maximum Edge). This procedure runs in linear time on the number
 227 of edges. In fact, the procedure does as many comparisons as the sum of vertex degrees which is two times the
 228 number of vertices. It uses a hash table m to store, for each vertex x , the maximal value of previously visited
 229 adjacent edges.

```

230 SLME procedure
231 Input : a valued graph (V,E,s)
232 Output : a relation R on V
233 L := {}
234 D := {}
235 for every x in V, m[x] := -1
236 while V-L is not empty
237   Select one vertex v in V-L
238   add v to L
239   C := {v}
240   while C is non empty
241     x := pop(C)
242     add x to L
243     add neighbours(x) - L to C
244     m[x] := max{s(n): n in neighbours(x)}
245     for every n in neighbours(x)
246       if m[n] = m[x] add {n,x} to R
  
```

250 Once done, the clustering phase consists in computing the reduced graph G/R , whose vertices are the con-
 251 nected components of the subgraph (V, R) of G and in inducing a new valuation according to a hierarchical
 252 criteria chosen among the following:

253 **single-link**: the value of an edge in G/R between two components C_1, C_2 is the maximal value of edges in
 254 $E_{C_1, C_2} = E \cap (C_1 \times C_2)$.

255 **complete-edge**: the minimal value in E_{C_1, C_2} ,

256 **average-edge**: the average value in E_{C_1, C_2} ,

257 **vertex-weight**: the sum of values in E_{C_1, C_2} over $|C_1| + |C_2|$

258

259 Observe that the above *complete-edge* and *average-edge* criteria differ from the usual complete and average
 260 link clustering since they are computed on a restricted set of pairs. The *vertex-weight* criterion is the one that
 261 best minimised the chain effect in our experiments. However in general, single link will also be satisfactory
 262 because the chain effect has already been reduced by the SLME procedure. In fact, this approach appears
 263 to be robust with regard to the clustering criteria. It is more sensitive to the existence of very small values

264 in the similarity matrix S . Indeed, any non null value will generate an edge in the graph and if this edge is the
 265 only one linking two objects, then these objects will be clustered even if the similarity is very small. This draw-
 266 back can be corrected by the use of a threshold which clarifies the borderline between null values and signifi-
 267 cant similarities.

268 The CPCL algorithm then becomes

```

269 Algorithm CPCL
270 input : a valued graph  $G = (V, E, s)$ 
271 parameters : a threshold  $t$  and a number of iterations  $I$ 
272 output : a partition of  $V$ 
273 for  $i = 1$  to  $i = I$  do
274    $E' := \{e \text{ in } E : v(e) > t\}$ 
275    $R := \text{SLME}(V, E', s)$ 
276    $G := G/R$ 
277 return  $V$ 
  
```

280 It involves I calls to the SLME procedure on the current reduced valued graph (V, E', s) .

281 It follows from this re-exploration of CPCL that it can be used for fast clustering of sparse similarity matrix
 282 with a reduced range of distinct values.

283 Until now, this algorithm has been applied to the following similarity matrix defined on groups of objects
 284 and generated in two steps:

285 **Step 1:** We consider a reduced subset of variation relations among those presented in Section 3.2 that we
 286 shall note COMP.

287 We then compute the set of connected components generated by the COMP relations. Terms that
 288 are not related by any of the variations in COMP will form singleton components.

289 **Step 2:** We select a second subset of variations denoted by CLAS to group components. Next, given two com-
 290 ponents C_1 and C_2 , a similarity value v is defined in the following way:

$$292 \quad v = \sum_{R \in \text{CLAS}} \frac{|R \cap C_1 \times C_2|}{|R|}$$

293 This similarity relies on the number of variations across the components and on the frequency of the variation
 294 type which on a large corpus will substantially reduce the influence of the most noisy variations like lexical
 295 substitutions on binary terms. The resulting matrix has all the characteristics that justify the application of
 296 the CPCL algorithm.

297

298 3.4. Implementation issues

299 Fig. 2 gives an overall view of the system. It is currently run on-line on a Linux Apache MySQL Php PERL
 300 Secured (LAMPS) server.⁴ The three components term extractor, relation identifier and clustering module are
 301 implemented as PERL5 OO programs while all the data are stored in a MySQL database. Clustering outputs
 302 can be accessed either via an integrated visualisation package (aiSee based on Graph Description Language)
 303 for domain topic mapping or through an interactive hypertext interface based on PERL DBI and CGI pack-
 304 ages. This interactive interface enables the user to browse the results, from the term network (variation links) to
 305 clusters contents and finally to documents where the terms appeared. The systems' modules can also be ex-
 306 cuted from this interface.

⁴ TermWatch is available for research purposes after obtaining an account and a password from the authors.

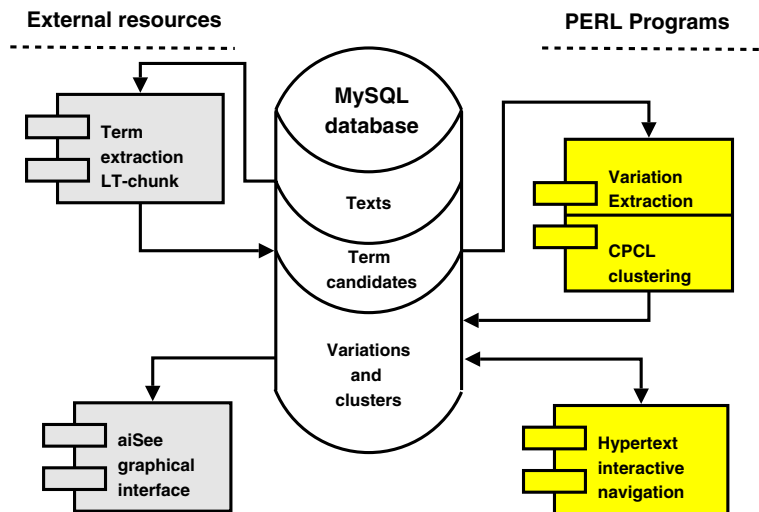


Fig. 2. Overall view of the TermWatch system.

307 4. Evaluation metrics

308 Evaluating the results of a clustering algorithm remains a bottleneck issue (Jain & Moreau, 1987; Tibsh-
 309 irani, Walther, & Hastie, 2000; Yeung & Ruzzo, 2001). The objective of the evaluation for our specific task:
 310 clustering multi-word terms out-of-context, is detailed in Section 4.1 followed by a review of existing evalua-
 311 tion methods Section 4.2. Finally, in Section 4.3 we propose enhancements to the editing distance suggested by
 312 Pantel and Lin (2002) for cluster evaluation.

313 4.1. Out-of-context Term Clustering (OTC)

314 Given a list of terms, the task consists in clustering them using exclusively surface lexical information in
 315 order to obtain coherent clusters. In this framework, clustering is done without contextual document informa-
 316 tion, without any training set and in a completely unsupervised way. We refer to this task as OTC (Out-of-
 317 context Term Clustering).

318 Let us emphasise that OTC is different from Entity Name Recognition (ENR). ENR task as described in
 319 Kim, Ohta, Tsuruoka, Tateisi, and Collier (2004) is based on massive learning techniques and new terms are
 320 forced to enter known categories. Whereas in unsupervised clustering, a new cluster can be formed of terms
 321 not belonging to an already existing category. This can lead to the discovery of new domain topics. It should
 322 also be noted that MWTs cannot be reduced to single words. Unlike single words, a MWT can occur only
 323 once “as is” (without variations) in the whole corpus. It is thus difficult for the usual *document* × *feature* rep-
 324 resentation to find enough frequency information to form clusters. Therefore methods based on *term-docu-*
 325 *ment* representation cannot be directly applied to OTC without adaptation. This adaptation is described in
 326 further details in Section 5.

327 4.2. Existing measures for cluster evaluation

328 Cluster evaluation generally falls under one of these two frameworks:

- 329 (1) Intrinsic evaluation: evaluation of the quality of the partitions vis-à-vis some criteria.
- 330 (2) Extrinsic evaluation: task-embedded evaluation or evaluation against a gold standard.

331

332 Intrinsic evaluation, also called “internal criteria” is used to measure the intrinsic quality of the clusters in
 333 the absence of an external ideal partition. Internal criteria concern measures like cluster homogeneity and

334 separation, or the stability of the partitions with respect to sub-sampling (Hur, Elisseeff, & Guyon, 2002).
 335 Alternatively, the measure can also seek to determine the optimal number of clusters (Hur et al., 2002).

336 Extrinsic evaluation, also known as “external criteria” refers to the comparison of a partition against an
 337 external ideal solution (gold standard) (Jain & Moreau, 1987; Milligan & Cooper, 1985) or a task-embedded
 338 evaluation. The comparison with a gold standard is done using measures like the Rand index or its adjusted
 339 variant (Hubert & Arabie, 1985) that measures the degree of agreement between two partitions.⁵ Milligan and
 340 Cooper (1986) recommended the use of Adjusted Rand index even when comparing clusters at different levels
 341 of the hierarchy. As observed by Yeung and Ruzzo (2001), external criteria has the advantage of providing
 342 an “independent unbiased assessment of the cluster” but has as inconvenience the fact that they are hardly
 343 available.

344 Internal criteria has as advantage the fact that it can bypass the necessity of having an external ideal solu-
 345 tion but its major inconvenience is that evaluation is based on the same information from which the clusters
 346 were derived. Pantel and Lin (2002) observed a flaw in the external criteria approach as suggested by the Rand
 347 index. According to them, computing the degree of agreements and disagreements between proposed parti-
 348 tions and an ideal one can lead to unintuitive results. For instance, if the ideal partition has 20 equally-sized
 349 clusters with 1000 elements each, treating each element as its own cluster will lead to a misleading high score of
 350 95%. We observe also that the Rand index and the adjusted Rand Index (Hubert & Arabie, 1985) have the
 351 following flaws:

- 352 • they are computationally expensive since they require $|\Omega|^2$ comparisons which is problematic when $|\Omega|$ is
 353 large,
- 354 • they are too sensitive to the number of clusters when comparing clustering outputs of different size (Weh-
 355 rens, Buydens, Fraley, & Raftery, 2003),
- 356 • the adjusted Rand Index supposes a hyper-geometric model which is obviously not fitted to the distribution
 357 of terms in the current experiment (GENIA categories).

358
 359 Denoeud, Garreta, and Guénoche (2005) tested the ability of different measures in determining the distance
 360 between two partitions. The Jaccard measure appeared as the best in this task since it does not have the draw-
 361 backs of the (adjusted) Rand Index. It computes the number of pair of items clustered together by two algo-
 362 rithms divided by the total number of pairs clustered by one of the algorithms. However, it cannot take into
 363 account the specificities of a target distribution. More precisely, suppose that we want to measure the gap
 364 between a clustering output and a target classification, suppose moreover that the target classification has a
 365 very large class with a great number of terms whereas the mean size of the other classes is small, (this is pre-
 366 cisely the case in the GENIA taxonomy where the *other name* class groups 33% of all the terms in this taxon-
 367 omy), although this class is disproportioned, it is definitely not the most informative. The Jaccard measure will
 368 favour methods that focus on the detection of the biggest class against more fine-grained measures that try first
 369 to fit the distribution of items in the smaller classes. Yeung and Ruzzo (2001) proposed a compromise for clus-
 370 ter evaluation in which evaluation is based on the predictive capacity of the methods vis-à-vis a hidden exper-
 371 imental condition. They tested their method on gene expression (microarray) data. This approach, aside from
 372 being computationally intensive, is not suitable for datasets where no experimental conditions (hidden or
 373 otherwise) obtain nor will it be suitable for datasets where the different samples do not share any dependent
 374 information.

375 In the task-embedded evaluation framework, what is evaluated is not the quality of the entire partition but
 376 rather that of the *best cluster* (Pantel & Lin, 2002), i.e., the cluster which enables the user to best accomplish
 377 his information seeking need. This is typically the case with cluster evaluation in the information retrieval
 378 field.

379 Following the extrinsic evaluation approach, Pantel and Lin (2002) proposed the use of the editing distance
 380 to evaluate clustering outputs. The idea is to evaluate the *cost* of producing the ideal solution from the pro-

⁵ Given two equivalence relations P and Q defined on a set Ω , the rand Index is the number of agreements between the two relations $|(P \cap Q) \cup \neg(P \cap R)|$ over the total number of pairs $|\Omega|^2$. The adjusted rand index assumes the generalized hypergeometric distribution as the model to ensure that two random partitions do take a constant null value.

381 posed partitions. This supposes the existence of an external ideal solution. The editing distance is an old notion
 382 used to calculate the cost of elementary actions like *copy*, *merge*, *move*, *delete* needed to obtain one word (or
 383 phrase or sentence) from another. Here, the authors applied it to cluster contents and chose to consider three
 384 elementary actions: copy, merge, move. Considering the OTC task, we needed a measure that focused on cluster
 385 quality (homogeneity) vis-à-vis an existing partition (here the GENIA classes). Pantel and Lin's editing
 386 distance appeared as the most suitable for this task. It is adapted to the comparison of methods producing
 387 a great number of clusters (hundreds or thousands) and of greatly differing sizes. On a more theoretical level,
 388 the idea of editing distance is conceptually suited to the nature of our evaluation task, i.e., calculate the *effort*
 389 or the *cost* required to attain an existing partition from the ones proposed by automatic clustering methods.

390 4.3. Metrics for evaluation of clusters

391 Given an existing target partition, Pantel and Lin's (2002) measure evaluates the ability of clustering algo-
 392 rithms to detect part of the structure represented by this partition. This measure extends the notion of editing
 393 distance to general families of subsets of items. In particular, it allows to consider fuzzy clustering where clus-
 394 ters overlap (copy action). Here we will not use this feature since we target crisp clustering. Hence, we focus on
 395 the two elementary operations: *merges* which is the union of disjoint sets and *moves* that apply to singular ele-
 396 ments. In this restricted context, Pantel and Lin's (2002) measure has a more deterministic behaviour and
 397 shows some inherent bias which we will correct.

398 To measure the distance between a clustering output and an ideal partition, these authors considered the
 399 minimal number of merges and moves that have to be applied to a clustering output in order to obtain the
 400 target partition. In fact, this number can be easily computed since the number of merges corresponds to
 401 the number of extra-classes and the number of moves to the number of elements that are not in the dominant
 402 class of the cluster. Indeed, each cluster is associated to the class with which it has the maximum intersection.
 403 The elements of a cluster which are not in the intersection will then have to be moved.

404 Thus, let Ω be a set of objects for which we know a crisp classification $\mathcal{C} \subseteq 2^\Omega$, seen as a family of subsets of
 405 Ω such that $\cup \mathcal{C} = \Omega$ and $C \cap C' = \emptyset$ for all C, C' in \mathcal{C} . Consider now a second disjoint family \mathcal{F} of subsets of
 406 Ω representing the output of a clustering algorithm. For each cluster $F \in \mathcal{F}$, we denote by \mathcal{C}_F the class $C \in \mathcal{C}$
 407 such that $|C \cap F|$ is maximal. Pantel and Lin's measure can be re-formulated thus:

$$410 \mu_{LP}(\mathcal{C}, \mathcal{F}) = 1 - \frac{(|\mathcal{F}| - |\mathcal{C}|) + \sum_{F \in \mathcal{F}} (|F| - |\mathcal{C}_F \cap F|)}{|\Omega|} \quad (1)$$

411 In the numerator of formula (1), the term $|\mathcal{F}| - |\mathcal{C}|$ gives the number *Mg* of necessary merges, and the sum
 412 $\sum_{F \in \mathcal{F}} (|F| - |\mathcal{C}_F \cap F|)$ the number *Mv* of moves. The denominator $|\Omega|$ of (1) is supposed to give the maximal
 413 cost of building the classification \mathcal{C} from scratch. Indeed, Pantel and Lin considered two trivial partitions: the
 414 discrete one where all clusters are singletons (every term is its own cluster) and the complete one where all
 415 terms are in a single cluster. These trivial partitions are supposed to be at equal distance from the target clas-
 416 sification. These authors suggest that the complete clustering needs $|\Omega|$ moves and the discrete $|\Omega|$ merges but
 417 this turns out not to be the case.

418 Clearly, discrete clustering only needs $|\Omega| - |\mathcal{C}|$ merges. Moreover, if $g = \max\{|C| : C \in \mathcal{C}\}$ is the size of the
 419 largest class in \mathcal{C} , then the distance of the trivial complete partition to the target partition is $|\Omega| - g$. It follows
 420 that in the case where g is much more greater than the mean size of classes in $|\mathcal{C}|$, Pantel and Lin's measure,
 421 based on the total number of necessary moves and merges over $|\Omega|$ favours the trivial complete partition over
 422 the discrete one and therefore algorithms that produce very few clusters, even of poor quality. Incidentally,
 423 this happens to be the case with the GENIA classes. Following these observations, we propose the following
 424 corrected version (2) where the weight of each move is no more 1 but $|\Omega|/(|\Omega| - g)$ and the weight of a merge is
 425 $|\Omega|/(|\Omega| - |\mathcal{C}|)$:

$$426 \mu_{ED}(\mathcal{C}, \mathcal{F}) = 1 - \frac{\max\{0, |\mathcal{F}| - |\mathcal{C}|\}}{|\Omega| - |\mathcal{C}|} - \frac{\sum_{F \in \mathcal{F}} (|F| - |\mathcal{C}_F \cap F|)}{|\Omega| - g} \quad (2)$$

$$428 = 1 - \frac{Mg}{|\Omega| - |\mathcal{C}|} - \frac{Mv}{|\Omega| - g} \quad (3)$$

429 The maximal value of μ_{ED} is 1 in the case where the clustering output corresponds exactly to the target par-
 430 tition. It is equal to 0 in the case that \mathcal{F} is a trivial partition (discrete or complete).

431 However, μ_{ED} can also take negative values. Indeed consider the extreme case where \mathcal{C} is of the form
 432 $\{A, B_1, \dots, B_n\}$ with one class $A = \{\alpha_1, \dots, \alpha_n, w_1, w_2\}$ with $n + 2$ elements and n singleton classes $B_i = \{\beta_i\}$.
 433 Now take as \mathcal{F} the whole family of n pairs $\{\alpha_i, \beta_i\}$ for $1 \leq i \leq n$ augmented with the singletons $\{w_1\}, \{w_2\}$.
 434 Then

$$436 \quad \mu_{ED}(\mathcal{C}, \mathcal{F}) = 1 - \frac{1}{(n + n + 2) - (n + 1)} - \frac{n}{(n + n + 2) - (n + 2)} = -\frac{n}{n + 1} < 0$$

437 and $\lim_{n \rightarrow \infty} \mu_{ED}(\mathcal{C}, \mathcal{F}) = -1$

438 In fact, in the case that g is much more greater than the mean size of classes and that the distribution of sizes
 439 of classes fits an exponential model, we have experimentally checked that $\mu_{ED}(\mathcal{C}, \mathcal{F}) \in]-1, 0[$ for random
 440 clusterings \mathcal{F} with $2g$ clusters and equiprobability for an element w to be affected to anyone of these clusters.

441 Based on the corrected μ_{ED} index, we propose a complementary index, Cluster homogeneity (μ_H) defined as
 442 the number of *savings* (product of μ_{ED} per $|\Omega|$) over the number Mv of movings:

$$444 \quad \mu_H(\mathcal{C}, \mathcal{F}) = \frac{\mu_{ED}}{1 + Mv} \times |\Omega|$$

445 μ_H takes its maximal value $|\Omega|$ if $\mathcal{F} = \mathcal{C}$ and, like the μ_{ED} measure, it is null if \mathcal{F} is one of the two trivial
 446 partitions.

447 We will use μ_H to distinguish between algorithms having similar editing distances but not producing clus-
 448 ters of the same quality (homogeneity). However, since the cluster homogeneity measure relies on the cor-
 449 rected editing distance (μ_{ED}), for a method to obtain a good cluster homogeneity measure (μ_H), it also has
 450 to show a good savings value (good μ_{ED}).

451 5. Experimental setup

452 In this section, we describe the principles (relations) used for clustering (Section 5.1), the different term rep-
 453 resentations adopted for the methods evaluated (Section 5.2) and the clustering parameters for each method
 454 (Section 5.3).

455 5.1. The relations used for clustering

456 Given the OTC task, our experiment consisted in searching for the principle and the method that can best
 457 perform this task. Three principles were tested:

458 **CLS:** Clustering by coarse lexical similarity: grouping terms simply by identical head word. We call this
 459 “baseline” clustering as it is technically the most straightforward to implement and is also a more basic
 460 relation than the ones used by TermWatch (see Section 3.2). However, it should be noted that this head
 461 relation is not so trivial for the GENIA corpus. Indeed, [Weeds, Dowdall, Keller, and Weir \(2005\)](#) showed
 462 that grouping terms by identical head words enables to form rather homogeneous clusters with regard to
 463 the GENIA taxonomy. In their experiment, out of 4, 797 clusters, 4104 (85%) contained terms with the
 464 same GENIA category while 558 (12%) clusters contained terms with 2 or 3 semantic categories. A further
 465 135 (3%) clusters contained terms with more than p semantic categories.

466 **LSS:** Clustering by fine-grained Lexico-Syntactic Similarity as implemented in the TermWatch system
 467 using the CPCL clustering algorithm described in Section 3.3. Terms are represented as a graph of
 468 variations.

469 **LC:** Clustering by Lexical Cohesion. This principle required a spatial representation based on a vector rep-
 470 resentation of terms in the space of words they contain. It was suggested by the characteristics of the base-
 471 line and graph (LSS) representations. The LC representation offers a numerical encoding of term similarity
 472 that allows us to subject statistical clustering approaches (hierarchical and partitioning algorithms) to the
 473 OTC task. We describe this representation in more details below.

476 5.2. Vector representation for statistical clustering methods

477 In order for statistical clustering methods to find sufficient *co-occurrence* information in an OTC task, it was
 478 necessary to represent *term-term* similarity. We redefined *co-occurrence* here as *intra-term* word *co-occurrence*
 479 and built a *term* \times *word* matrix where the rows were the terms and the columns the unique constituent words.

480 To ensure that the statistical methods will be clustering on a principle as close as possible to the LSS rela-
 481 tions used by TermWatch and to the head relation used by the baseline, we further adapted this matrix as fol-
 482 lows: words were assigned a weight according to their grammatical role in the term and their position with
 483 regard to the head word. Since a head word is the noun focus (the subject), it receives a weight of 1. Modifier
 484 words are assigned a weight which is the inverse of their position with regard to the head word. For instance,
 485 given the term “*coronary heart disease*”, *disease* (the head word) will receive a weight of 1, heart will be
 486 weighted 1/2 and coronary 1/3.

487 More formally, let $W = (w_1, \dots, w_N)$ be the ordered list of words occurring in the terms. A term
 488 $t = (t_1, \dots, t_q)$ can be simply viewed (modulo permutations) as a list of words where the t_i are words, t_q
 489 the head and t_1, \dots, t_{q-1} is a possible empty list of modifiers. Each term t is then associated with the vector
 490 V_t such that

$$492 \quad V_t[i] = \begin{cases} \frac{1}{1+q-j} & \text{whenever } w_i = t_j \\ 0 & \text{elsewhere} \end{cases}$$

493 Let M be the matrix whose rows are the V_t vectors. We derive two other matrices from M :

- 494 (1) A similarity matrix $S = M \cdot M^t$ whose cells give the similarity between two terms as the scalar product of
 495 their vectors (for hierarchical algorithms).
 496 (2) A core matrix C by removing all rows of M corresponding to terms with less than three words and all
 497 columns corresponding to words that appeared in less than 5% of the terms. Indeed, experimental runs
 498 showed that the k -means algorithms could not produce meaningful clusters when considering the matrix
 499 of all terms.

500
 501 This weighting scheme translates the linguistic intuition that the further a modifier word is from the head,
 502 the weaker the semantic link with the concept represented by the head. This idea shares some fundamental
 503 properties with the relations used by TermWatch for clustering. Note also that this weighting scheme is a more
 504 fine-grained principle than the one used by the baseline. Representing terms in this way leads to the identifi-
 505 cation of lexically-cohesive terms (i.e., terms that often share the same words). This idea was explored by
 506 Dobrynin, Patterson, and Rooney (2004) although in a different way. Their *contextual document clustering*
 507 method focused on the identification of words that formed *clusters of narrow scope*, i.e. lexically cohesive
 508 words which appeared with only a few other words. Lexical cohesion is not a new notion in itself. It has
 509 already been explored in NLP applications for extracting collocations (fixed expressions) from texts (Church
 510 & Hanks, 1990; Smadja, 1993).

511 5.3. Clustering parameters

512 MWTs were clustered following the three types of relations described in Section 5.1. The following methods
 513 were tested: baseline; CPCL on graph of variations; partitioning (k -means, Clara based on medoids), hierar-
 514 chical (CPCL on similarity matrix S).

- 515 • **Baseline on CLS:** No particular parameter is necessary. All terms sharing the same head word are put in the
 516 same cluster.
 517 • **CPCL on LSS:** Parameter setting consists in assigning a role to each relation (*COMP* or *CLAS*). Among all
 518 the variations extracted by TermWatch, we selected a subset that optimised the number of terms over
 519 the maximal size of a class. Hence this selection was done without prior knowledge of the GENIA taxon-

omy. The variations selected for the *COMP* phase are those where terms share the same head word or WordNet semantic variants. In the current experiment, by order of ascending cardinality, *COMP* relations were:

- spelling variants,
- substitutions of modifiers filtered out using WordNet (*sub_wn_modifier*),
- insertion of one modifier word (*strong_ins*),
- addition of one modifier word to the left (*strong_exp_1*),
- substitutions of the first modifier in terms of length ≥ 3 (*strong_sub_modifier_3*).

The *CLAS* variations were

- WordNet head substitutions (*sub_head_wn*),
- insertions of more than one modifier (*weak_ins*),
- addition of more than one modifier word to the left (*weak_exp_1*),
- substitution of modifiers in terms of length ≥ 3 (*weak_sub_modifier_3*).

No threshold was set so as not to exclude terms and relations. Since the objective of this experiment is to form clusters as close as possible to the GENIA classes, the algorithm was stopped at iteration 1. Thus, only a few part of relations induced by the variations were really used in the clustering. More precisely, only relations induced by rare variations which are assigned a higher weight or relations between near-isolated terms were considered. Hence, the exact technique used in agglomerative clustering (single, average or complete link) did not come into play here. We also tested the performance of the 1st step grouping, i.e., the level forming connected components (*COMP*) with a subset of the relations. This level is akin to baseline clustering although the relations are more fine-grained.

• **Hierarchical on LC:** Clustering is based on the similarity matrix $S[S \geq th]$ derived from S by setting to 0 all values under a threshold th . We used the following values for th :

- 0.5: the rationale is that at this weight, terms either share the same head or have common modifiers close to the head,
- 0.8: this weight imposes the same head on related terms.

Because the dissimilarity matrix was too large, we had to use our own PERL programs to handle such sparse matrices. Based on a graph representation of the data, only non zero values were stored as edge values enabling each iteration to be done in a single search. We were thus able to run the usual variants of single, average and complete link hierarchical clustering on this system but they did not produce any relevant clustering (all the cluster evaluation measures were negative). Since the similarity matrix S had all the requirements to be an input to the CPCL algorithm, we subjected it to the CPCL algorithm. After some tests, we finally selected the *vertex-weight* (Section 3.3) as the agglomerative criterion since it significantly reduced the chain effect. We did four iterations for each threshold value. This yielded significant results. Thus the results shown for hierarchical clustering were obtained using the CPCL algorithm on the *term* \times *word* matrix.

• **Partitioning on LC:** This method is based on the computation of k -means centers and medoids on the core matrix C . We used the standard functions of k -means and CLARA (Clustering LARge Applications) fully described in Kaufman and Rousseeuw (1990). CLARA considers samples of datasets of fixed size on which it finds k medoids using PAM algorithm (Partitioning Around Medoids) and selects the results that induce the best partition on the whole dataset. PAM is supposed to be a more robust version of k -means because it minimizes a sum of dissimilarities instead of a set of distances. However, for large datasets, PAM cannot be directly applied since it requires a lot of computation time. CLARA and PAM are available on the standard R cluster package.⁶ To initialize CLARA, we used the same procedure as CLARANS (Ng & Han, 2002) to draw random samples using PERL programs and a graph data structure. We ran these two variants (k -means and CLARA) for the following values of k : 36, 100, 300, 600 and 900. Then, given these centers and medoids, we again used our PERL programs for storing large sparse matrix, to assign each term to its nearest center or medoid and to obtain a partition on the whole set of terms.

⁶ Version 1.10.2, 2005-08-31, by Martin Maechler, based on S original by Peter Rousseeuw (rousse@uia.ua.ac.be), Anja.Struyf@uia.ua.ac.be and Mia.Hubert@uia.ua.ac.be.

568

569 The results of clustering with these algorithms and their variants were then evaluated against the target par-
 570 tition (the GENIA taxonomy) using the measures described in Section 4.3. Combining *R* and PERL 5 has
 571 been quite efficient. *R* offers very robust implementations of spatial clustering algorithms while PERL allows
 572 one to easily define optimal data structures. Thus all the data processing including the initialization phase and
 573 sample extraction was done with PERL, leaving to *R* the massive numerical computations based on C and
 574 FORTRAN subroutines. All the tests were performed on a PENTIUM IV PC server running LINUX
 575 DEBIAN stable with 1Go of RAM, SCSI disk and no X11 server for memory saving.

576 6. Results

577 6.1. Possible impact of the variations on TermWatch's performance

578 Before comparing the clustering results obtained by the different methods, we investigated the possible
 579 impact of the variations used by TermWatch on its performance. The idea was to determine if our variation
 580 relations alone could reproduce these categories, i.e., if they grouped together terms from one only GENIA
 581 class. In this case, then there would be no need to perform clustering since the variation relations alone
 582 can discover the ideal partition. However, our study showed that this was not the case.

583 The chart Fig. 3 shows for each of our variation relation, the number of links acquired, the proportion of
 584 intra-category links and the proportion inter-category links (from different classes). We can see clearly from
 585 this figure that some relations are rare, i.e., they capture too few links although they link terms from the same
 586 class (*sub_modifier_wn*, *strong_ins*, *weak_ins*). These relations are in the minority especially by the proportion of
 587 terms linked. Other relations like *weak_exp2*, *weak_sub_head3*, *weak_exp_r* are more abundant but they lead
 588 to heterogeneous clusters, they link terms from different GENIA classes. Surprisingly, *weak_exp_l* and
 589 *strong_sub_mod3* produced relatively good quality clusters while relating a considerable number of terms.

590 6.2. Evaluation of clustering results

591 Using the relations chosen in Section 5.3, CPCL on LSS generated 1897 non trivial components (at the
 592 COMP phase) involving only 6555 terms. Adding CLAS relations in the second phase led to 3738 clusters
 593 involving 19,887 terms.

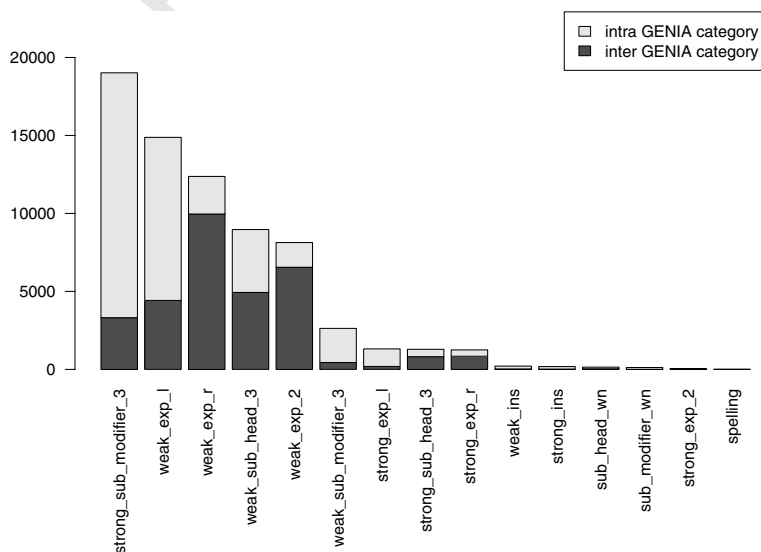


Fig. 3. Distribution of related pairs of terms by variations.

594 Hierarchical clustering based on similarity matrix introduced in Section 5.2 generated 1090 clusters involv-
 595 ing 25.129 terms for a threshold $th = 0.5$ and 1217 clusters involving 19,867 terms for $th = 0.8$.

596 The plots in Figs. 4 and 5 show the results of the evaluation measures μ_{ED} and μ_H introduced in Section 4.3.
 597 Since the majority of the clustering methods are sensitive to term length, we plotted the score obtained by
 598 each of the measure (y-axis) by term length (x-axis). Note that at each length, only terms of that length
 599 and above are considered. For instance, at length 1, all terms are considered. At length 2, only terms having
 600 at least two words are considered. Thus, the further we move down the x-axis, the fewer the input terms for
 601 clustering.

602 Fig. 4 shows the % of savings obtained by the nine algorithms tested using the corrected ED measure. We
 603 see that the hierarchical method with a threshold = 0.8 and CPCL obtain a better score than the baseline clus-
 604 tering when considering all the terms (length ≥ 1). When fewer and longer terms are considered (length ≥ 3),
 605 partitioning methods outperform CPCL and hierarchical algorithms but still remain below the baseline. This
 606 is because, at length ≥ 3 , CPCL has fewer terms, thus fewer relations with which to perform the clustering.
 607 Statistical methods on the other hand, with longer terms have a better context, thus more relations in the
 608 matrix. From terms of length ≥ 4 words, partitioning methods outperform the baseline.

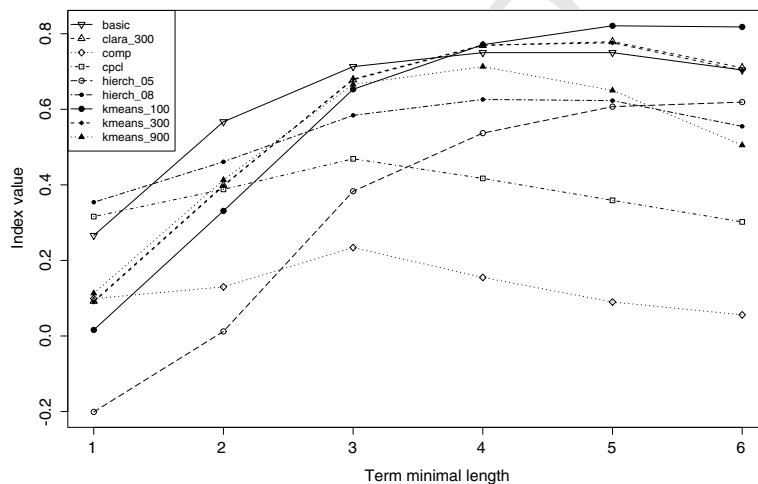


Fig. 4. Editing distance between clustering results μ_{ED} and Genia categories.

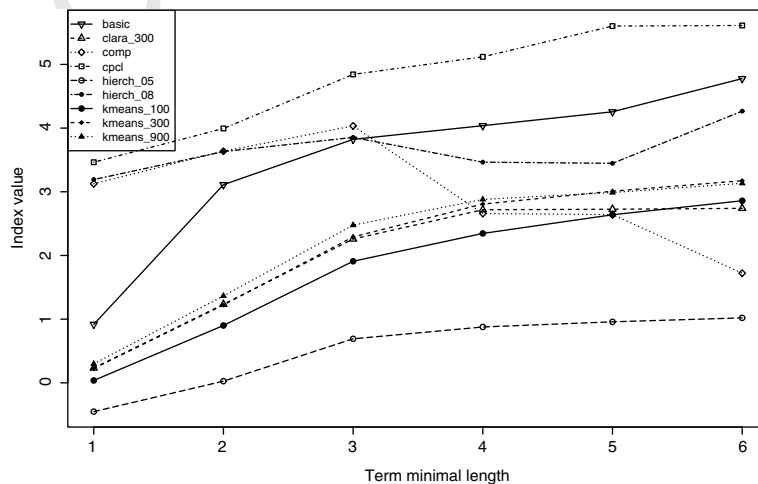


Fig. 5. Cluster homogeneity measure μ_{ED} on the Genia categories.

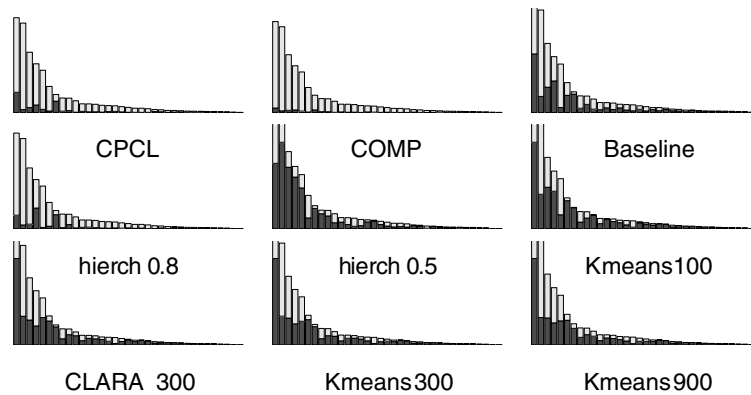


Fig. 6. Proportion of intra- and inter-GENIA category terms per algorithm. The black bars represent mis-classifications. White bars represent terms from the same GENIA category.

609 However, the ED measure masks important features of the clustering outputs since it is a compromise
 610 between the number of necessary moves and merges needed to reach the target partition. More important
 611 is the quality of the clusters (cluster homogeneity) vis-à-vis the target partition (GENIA classes). This is mea-
 612 sured by the μ_H which calculates the ratio between the value of ED and the number of movings. The μ_H per-
 613 formance of the algorithms is shown in the plot of Fig. 5.

614 It appears clearly that on cluster quality, CPCL is the only algorithm that significantly outperforms the
 615 baseline irrespective of term length. Hierarchical algorithm with $th = 0.8$ and the COMP phase of CPCL fol-
 616 low closely but only on all terms (length ≥ 1). Their performance drops when terms of length ≥ 4 are con-
 617 sidered. Partitioning algorithms show poor cluster homogeneity. *K*-means with $k = 100$ performs worse than
 618 the other variants (Clara, *k*-means with $k = 300$, $k = 900$). Hierarchical with $th = 0.5$ obtain the poorest score.

619 To gain a better insight on the cluster homogeneity property, we generated for every algorithm a chart
 620 showing the proportion of terms which share the same GENIA class with the majority of terms in the same
 621 cluster (and thus that do not require any move) The nine charts are shown in Fig. 6.

622 It appears that the COMP variant of CPCL produced the most homogeneous clusters which is not alto-
 623 gether surprising because the relations used in COMP phase are the most semantically tight. COMP and
 624 CPCL significantly outperform the baseline. This good performance is a bit unexpected for CPCL because
 625 the CLAS relations induce a change of head word which could lead to a semantic gap (change of semantic
 626 class).

627 Closely following is the hierarchical algorithm at $th = 0.8$. The baseline comes fourth which shows that
 628 grouping terms simply by identical head words as done by baseline is good but not good enough to form
 629 semantically homogeneous clusters.

630 Partitioning methods clearly produced less homogeneous clusters. These algorithms showed low error rates
 631 roughly on categories with a low proportion of one word terms.

632 7. Concluding remarks

633 We have developed an efficient text mining system based on meaningful linguistic relations which works
 634 well on MWTs and thus on very large and sparse matrices. This method is suitable for highlighting rare phe-
 635 nomena which may correspond to weak signals.

636 The specific evaluation framework set up here led us to redefine a matrix representation in order to enable
 637 comparison with existing statistical methods. We defined a new term weighting scheme in the matrix represen-
 638 tation enabling statistical methods to build significant clusters. We also corrected an existing cluster evaluation
 639 measure and defined a complementary one focused on cluster homogeneity.

640 The choice of the evaluation metric made it possible to compare algorithms outputting very high number of
 641 clusters, with considerable differences in this number (between 100 for *K*-means and 3738 for CPCL). This was

642 done without any assumption of equal cluster size. We believe these differences did not handicap any algo-
 643 rithm unduly since all produced clusters whose numbers were very far from the target partition (36 classes),
 644 especially our own method. As we cannot define a priori the number of optimal clusters, CPCL's performance
 645 was hampered for the μ_{ED} measure. Statistical methods (both hierarchical and partitioning) were more sensi-
 646 tive to term length.

647 The results however show that CPCL performs well in terms of cluster quality (homogeneity). Since this
 648 approach is computationally tractable in linear time, it also appears to be the best candidate for tasks requir-
 649 ing interaction with users in real time, like interactive query refinement. This aspect will be explored in a sep-
 650 arate study.

651 Overall, this experiment has shown that even without adequate context (document co-occurrence), cluster-
 652 ing algorithms can be adapted to partially reflect a human semantic categorisation of scientific terms.

653 Another interesting finding of this study is that when considering an OTC or a similar task, it may be inter-
 654 esting to first consider clustering by a basic relation before resorting to more complex and fine-grained term
 655 representation. The performance of the baseline clustering in our experiment is far from poor. It could be sat-
 656 isfactory for some tasks, for instance as a first stage for learning new taxonomy or knowledge structures from
 657 texts. These can be further refined using more sophisticated approaches: fine-grained linguistic relations,
 658 machine learning techniques with manually tagged learning sets.

659 Appendix A. Example of rule used in term extraction

660 This following simple rule translates the hypothesis that the preposition “of” plays a major role in the for-
 661 mation of terms. Thus prepositional phrases introduced by this preposition are attached to their governing
 662 NP.

663 From the tagged sentence:

```
664 [[The_DT inability_NN ]] of_IN [[ E1A_CD gene_NN products_NNS]]
665 ( to_TO induce_VB )
666 [[ cytolytic_JJ susceptibility_NN ]].
667
```

668 Our term extraction module would extract two multi-word terms (MWTs):

```
669 [[ The_DT inability_NN ]] of_IN [[ E1A_CD gene_NN products_NNS]]
670 [[ cytolytic\_JJ susceptibility\_NN ]]
```

671 This rule can be formulated as the following regular expression:

```
673 If:
674 <mod>* <N>+ of <mod>* <N>+ <prepl> <verb> <mod>* <N>+
675 then return:
676 (1) <mod>* <N>+ of <mod>* <N>+
677 (2) <mod>* <N>+
678 where:
679 <mod> = a determiner (DT) and/or an adjective (JJ)
680 <N> = any of the noun tags (NN, NNS, NNPS, NNP)
681 <prepl> = all the prepositions excluding ‘of’
682 * = Kleene’s operator (zero or n occurrences of an item)
683 + = at least one occurrence
684
```

685 Appendix B. Variation rules

686 For the sake of clarity, all the variation rules will be given for the compound structure only.

687 *B.1. Lexical variants*688 **Modifier substitutions (sub_modifier)** can be identified with this simple rule:

689 t2 is a substitution of t1 if and only if:

690 $t1 = M m M' h$ and $t2 = M m' M' h$ 691 with $m' \langle \rangle m$

692 where

693 t1 and t2 are terms,

694 M and M' are optional sequence of modifier words,

695 m and m' are a modifier words

696 h and h' are head word.

697

698 A chain of modifier substitutions can highlight properties around the same concept. For instance, the fol-
699 lowing variants all specify a type of *human cell line*:700 • *human leukemia cell line*701 • *human lymphoblastoid cell line*702 • *human monoblastic cell line*703 • *human monocytic cell line*

704

705 **Head substitutions (sub_head)** are identified via the following rule:

706 t2 is a substitution of t1 if and only if:

707 $t1 = M m h$ and $t2 = M m h'$ 708 with $h' \langle \rangle h$

709

710 Head substitutions highlight on the other hand families of concepts sharing the same property:

711 • *tumor cell killling*712 • *tumor cell line*713 • *tumor cell nuclei*714 • *tumor cell proliferation*715 • *tumor cell type*

716

717

718 *B.2. Syntactic variants*

719 These rules identify the three types of expansion variants.

720 • **Left expansion (exp_l)**

721 t2 is a left-expansion of t1 if and only if:

722 $t1 = M h$ and $t2 = M' m' M h$ 723 For example, *Ad2 infection* has as left expansion *adenovirus 2 (Ad2) infection*.724 • **Insertion (ins)**

725 t2 is an insertion of t1 if and only if:

726 $t1 = M1 m M2 h$ 727 $t2 = M1 m m' M' M2 h$ 728 For instance, *CD3-stimulated T lymphocyte* has as insertion variant, the term *CD3-stimulated human periph-*
729 *eral T lymphocyte*. Modifier expansions enable us to create *generic – specific* links. Head expansions identify
730 topical shifts as in *human disease* and *human disease syndrome*. Equivalent terms which undergo either a

731 syntactic transformation like permutation variants (*information retrieval* ↔ *retrieval of information*) are also
 732 identified.

733 • *Right expansion* (exp_r)

734 t2 is a right-expansion of t1 if and only if:

$$735 \quad t1 = M h \text{ and } t2 = M h M' h'$$

736 An example of right expansion would be *B cell development* and *B-cell development and differentiation*. Left
 737 and right expansions (exp-2) can be combined in the same term to yield *left-right expansion*. An example
 738 would be the link between *AIDS* and *second AIDS retrovirus*.

739

740 **References**

- 741 Berry, A., Kaba, B., Nadif, M., SanJuan, E., & Sigayret, A. (2004). Classification et désarticulation de graphes de termes. In *Proceedings of*
 742 *the 7th international conference on textual data statistical analysis (JADT 2004)*, Louvain-la-Neuve, Belgium, pp. 160–170.
- 743 Braam, R., & Moed, H. A. A. V. R. (1991). Mapping science by combined co-citation and word analysis. 2. Dynamical aspects. *Journal of*
 744 *the American Society for Information Science*, 42(2), 252–266.
- 745 Callon, T., Courtial, J., & Laville, F. (1991). Co-word analysis as a tool for describing the network of interactions between basic and
 746 technological research: the case of polymer chemistry. *Scientometrics*, 22(1), 155–205.
- 747 Church, K. W., & Hanks, P. (1990). Word association norms, mutual information and lexicography. *Computational Linguistics*, 16(1),
 748 22–29.
- 749 Cutting, D., Karger, D., Pedersen, J., & Tukey, O. (1992). Scatter/Gather: a cluster-based approach to browsing large document
 750 collections. In *15th annual international conference of ACM on research and development in information retrieval—ACM SIGIR*,
 751 *Copenhagen, Denmark*, pp. 318–329.
- 752 Denoeud, L., Garreta, H., & Guénoche, A. (2005). Comparison of distance indices between partitions. In P.L. et al. (Ed.), *Proceedings of*
 753 *applied stochastic models and data analysis*, Brest, pp. 17–20.
- 754 Dobrynin, V., Patterson, D., & Rooney, D. (2004). Contextual document clustering. In *Proceedings of the European conference on*
 755 *information retrieval (ECIR'04)*, Sunderland, UK, pp. 167–180.
- 756 Dunning, T. (1993). Accurate methods for statistics of surprise and coincidence. *Computational Linguistics* (19), 61–74.
- 757 Eisen, M., Spellman, P., Brown, P., & Botstein, D. (1998). Cluster analysis and display of genome-wide expression patterns. *Proceedings of*
 758 *the National Academy of Science, USA*(95), 14863–14868.
- 759 Fellbaum, C. (Ed.). (1998). *WordNet, an electronic lexical database*. MIT Press.
- 760 Glenisson, P., Glänzel, W., Janssens, F., & Moor, B. D. (2005). Combining full text and bibliometric information in mapping scientific
 761 disciplines. *Information Processing and Management*, 41(6), 1548–1572.
- 762 Hubert, L., & Arabie, P. (1985). Comparing partitions. *Journal of Classification*, 193–218.
- 763 Hur, B., Elisseeff, A., & Guyon, I. (2002). A stability-based method for discovering structure in clustered data. *Pacific Symposium on*
 764 *Biocomputing* (7), 6–17.
- 765 Ibekwe-SanJuan, F. (1998a). A linguistic and mathematical method for mapping thematic trends from texts. In *Proceedings of the 13th*
 766 *European conference on artificial intelligence (ECAI)*, Brighton, UK, pp. 170–174.
- 767 Ibekwe-SanJuan, F. (1998b). Terminological variation, a means of identifying research topics from texts. In *Proceedings of joint ACL-*
 768 *COLING'98, Québec*, pp. 564–570.
- 769 Ibekwe-SanJuan, F., & SanJuan, E. (2004). Mining textual data through term variant clustering: the termwatch system. In *Proceedings of*
 770 *recherche d'Information assistée par ordinateur (RIAO)*, Avignon, pp. 26–28.
- 771 Jacquemin, C. (2001). *Spotting and discovering terms through Natural Language Processing*. MIT Press.
- 772 Jain, A., & Moreau, J. (1987). Bootstrap technique in cluster analysis. *Pattern Recognition*, 20, 547–568.
- 773 Karypis, G., Han, E., & Kumar, V. (1994). Chameleon: a hierarchical clustering algorithm using dynamic modeling. *IEEE Computer:*
 774 *Special issue on Data Analysis and Mining*, 32(8), 68–75.
- 775 Kaufman, L., & Rousseeuw, P. (1990). *Finding groups in data: An introduction to cluster analysis*. John Wiley & Sons.
- 776 Kim, J.-D., Ohta, T., Tsuruoka, Y., Tateisi, Y., & Collier, N. (2004). Introduction to the bio-entity recognition task at JNLPBA. In *The*
 777 *proceedings of JNLPBA-04*, pp. 70–75.
- 778 Milligan, G. W., & Cooper, M. (1985). An examination of procedures for determining the number of clusters in a data set. *Psychometrika*,
 779 50, 159–179.
- 780 Milligan, G. W., & Cooper, M. (1986). A study of the comparability of external criteria for hierarchical cluster analysis. *Multivariate*
 781 *Behavioural Research*, 21, 441–458.
- 782 Ng, R., & Han, J. (2002). Clarans: a method for clustering objects or spatial data mining. In *IEEE transactions on knowledge and data*
 783 *engineering*, vol. 14.
- 784 Pantel, P., & Lin, D. (2002). Clustering by committee. In *Annual international conference of ACM on research and development in*
 785 *information retrieval—ACM SIGIR, Tampere, Finland*, pp. 199–206.
- 786 Polanco, X., Grivel, L., & Royauté, J. (1995). How to do things with terms in informetrics: terminological variation and stabilization as
 787 science watch indicators. In *Proceedings of the 5th international conference of the international society for scientometrics and*
 788 *informetrics, Illinois, USA*, pp. 435–444.

- 789 Price, L., & Thelwall, M. (2005). The clustering power of low frequency words in academic webs. *Journal of the American Society for*
790 *Information Science and Technology*, 56(8), 883–888.
- 791 Sanjuan, E., Dowdall, J., Ibekwe-Sanjuan, F., & Rinaldi, F. (2005). A symbolic approach to automatic multiword term structuring.
792 *Computer Speech Language (CSL)*, 19(4), 524–542.
- 793 Smadja, F. (1993). Retrieving collocations from text: Xtract. *Computational Linguistics* (19), 143–177.
- 794 Tibshirani, R., Walther, G., & Hastie, T. (2000). Estimating the number of clusters in a dataset via the gap statistic. Technical Report. No.
795 208. Dept. of Statistics, Stanford University.
- 796 Weeds, J., Dowdall, J., Keller, G. S., & Weir, D. (2005). Using distributional similarity to organise biomedical terminology. *Terminology:*
797 *Special Issue on Application-driven Terminology Engineering*, 11(1), 107–141.
- 798 Wehrens, R., Buydens, L. M., Fraley, C., & Raftery, A. E. (2003). Model-based clustering for image segmentation and large datasets via
799 sampling. Tech. Rep. 424, Department of Statistics, University of Washington.
- 800 Yeung, K., & Ruzzo, W. (2001). Details of the adjusted rand index and clustering algorithms. supplement to the paper “an experimental
801 study on principal component analysis for clustering gene expression data”. *Bioinformatics* (17), 763–774.
- 802 Zitt, M., & Bassecouard, E. (1994). Development of a method for detection and trend analysis of research fronts built by lexical or co-
803 citation analysis. *Scientometrics*, 30(1), 333–351.
- 804