# Using Syntactic Contexts for Measuring Word Similarity

**Caroline Gasperin⋆ Pablo Gamallo† Alexandre Agustini† Gabriel Lopes† Vera de Lima⋆**

⋆Faculdade de Informática
PPGCC, PUCRS
Av. Ipiranga, 6681
90619-900 Porto Alegre, RS, Brasil
{caroline,vera}@inf.pucrs.br

†Faculdade de Ciências e Tecnologia
Universidade Nova de Lisboa
Departamento de Informática, Quinta da Torre
P-2825-114 Monte da Caparica, Portugal
{gamallo,aagustini,gpl}@di.fct.unl.pt

## Abstract

This paper explores different strategies for extracting similarity relations between words from parsed text corpora. The strategies we have analysed do not require supervised training nor semantic information available from general lexical resources. They differ in the amount and the quality of the syntactic contexts used to compare words. The paper presents in details the notion of syntactic context and how syntactic information could be used to extract semantic regularities of word sequences. Finally, experimental tests with a Brazilian Portuguese corpus demonstrate that similarity measures based on fine-grained and elaborate syntactic contexts perform better than those based on poorly defined contexts.

## 1 Introduction

The strategies for extracting semantic information from corpora can be roughly divided into two categories, knowledge-rich and knowledge-poor methods, according to the amount of knowledge they presuppose (8). Knowledge-rich approaches require some sort of previously encoded semantic information (9; 6; 3): domain-dependent knowledge structures, semantic tagged training copora, and/or semantic resources such as handcrafted thesauri: Roget's thesaurus, WordNet, and so on. Therefore, knowledge-rich approaches may inherit the main shortcomings and limitations of man-made lexical resources: limited vocabulary size, since they can include unnecessary general words, or do not include necessary domain-specific ones; unclear classification criteria, since their word classification is sometimes too coarse and does not provide sufficient distinction between words, or is sometimes unnecessarily detailed; and, obviously, considerable time and effort required by building thesauri by hand. By contrast, knowledge-poor approaches use no presupposed semantic knowledge for automatically extracting semantic information. These techniques can be characterised

as follows: no domain- specific information is available, no semantic tagging is used, and no static sources as dictionaries or thesauri are required. They attempt to extract the frequency of co-occurrences of words within various contexts to compute semantic similarity among words. More precisely, the similarity measure takes into account the contexts that words share or do not share, as well as the importance of these contexts for each word. Words which share a great number of contexts are considered as similar.

Since contexts can be defined in two different ways, two specific knowledge-poor strategies can also be distinguished: windows-based and syntactic-based techniques. Windows-based techniques consider an arbitrary number of words around a given word as forming its window, i.e., its context. The linguistic information about part-of-speech categories and syntactic groupings is not taken into account to characterise word contexts (1; 10). The syntactic-based strategy, on the contrary, requires specific linguistic information to specify the word context. First, it requires a part-of-speech tagger for assigning a morphosyntactic category to each word of the corpus. Then, the tagged corpus is segmented into a sequence of basic phrasal groupings (or chunks). Finally, simple attachment heuristics are used to specify the relations between and within the phrasal groupings. Once the syntactic analysis of the corpus is reached, each word in the corpus is associated to a set of syntactic contexts. Then, a statistical method compares the frequency of the shared contexts to judge word similarity (7; 4; 2). In both strategies, window-based and syntactic-based techniques, words will be compared to each other in terms of their contextual distribution; yet, we consider that syntactic analysis opens up a much wider range of more precise contexts than does simple windows strategy. As syntactic contexts represent linguistic dependencies involving specific semantic relationships, they should be considered as fine-grained clues for identifying semantically related words.

Since syntactic contexts can be defined in different ways, syntactic-based approaches can also be significantly different. Different pieces of linguistic informa-

tion can be taken into account to characterise syntactic contexts[1]. For instance, the information used by Lin (4) to define the notion of syntactic context is not the same than that used by Grefenstette (7). Nevertheless, the choice of a particular type of syntactic context for measuring word similarity has not been properly justified by those researchers.

This way, the main objective of this paper is to analyse the appropriateness or the inadequacy of different types of syntactic contexts for computing word similarity. More precisely, various syntactic-based strategies will be compared on the basis of different definitions of the notion of syntactic context.

For this purpose, we apply these strategies on a Brazilian Portuguese corpus from NILC (Interinstitutional Center of Computational Linguistics - USP/São Carlos/Brazil), with news documents. Experiments concerning different syntactic relations repported below illustrate and show significant differences among results.

The article is organised as follows. In the next section, various types of syntactic contexts will be analised. Special attention will be paid for the notion of syntactic context used by Grefenstette, as well as for the specific notion that we have defined. Then, in section 3.1, we will use the same statistical similarity measure to compare the appropriateness of the syntactic contexts defined in the previous section. The best results are obtained when the syntactic-based strategy relies on our notion of syntactic context. Samples of the results we have obtained are presented in the Appendix.

## 2 Types of Syntactic Contexts

In this section, we analyse the notion of syntactic context used by Grefenstette to compute similar words (7). Then, we extract further syntactic information from the partially parsed text in order to make syntactic contexts more elaborate. As a consequence, we obtain fine-grained contexts which contain more specific information than the one provided by Grefenstette's approach.

### 2.1 The Notion of Attribute by Grefenstette

Grefenstette calls "attributes" the syntactic contexts of a word. Attributes are extracted from binary syntactic dependencies between two words within a noun phrase or between the noun head and the verb head of two related phrases. A binary syntactic dependency could be noted:

$$< R, w1, w2 >$$

where R denotes the syntactic relation itself (e.g., ADJ, NN, NNPREP, SUBJ, DOBJ, and IOBJ), and w1 and w2 represent two syntactically related words. Table 1 shows some syntactic dependencies between the noun "cause" and other related words.

Then, for each word found in the text, the system selects the words that are syntactically related to it. The syntactically related words are considered the attributes of the given word, i.e., its syntactic contexts. For instance, a noun can be syntactically related to an adjective by means of the ADJ relation, to another noun by means of the NN and NNPREP relations, or to a verb by means of SUBJ, DOBJ, and IOBJ relations. These related words are taken to be the known attributes of the noun.

In order to select the attributes of "cause", the system takes as input all the binary dependencies between "cause" and other words. Then, it extracts all the specific words syntactically related to "cause", since they represent its particular attributes. For example, from the 4 dependencies illustrated in 1 between "cause" and another word, it is possible to extract 4 attributes of "cause" (see table 2).

In the Grefenstette's notation, the attributes extracted from noun modifiers (namely NN, ADJ, and NN-PREP modifiers) do not keep the name of the particular syntactic relation. So, *<jaundice>*, *<possible>*, and *<death>* are attributes of "cause" even though the syntactic relations NNPREP, ADJ and NN are not explicitly represented. When extracting verbal complements, though, the specific syntactic relation is still available: *<DOBJ, determine>* is a verbal attribute constituted by both the word related to "cause" (i.e. the verb "determine") and the specific syntactic relation DOBJ.

### 2.2 Underspecified Attributes

The notion of attribute defined in the previous section does not inherit all the available syntactic information from binary dependencies. Consider one of the Portuguese expressions found in our corpus: "autorização à empresa" (*permission to the company*). From this expression, *<empresa>* (*company*) is extracted as the attribute of "autorização" (*permission*). Yet, relevant information implicitly contained in the dependency relation has been lost:

- information about the specific preposition: the attribute *<empresa>* does not convey information about the particular preposition "a" relating the two words;

- information about the opposite attribute: the attribute *<autorização>* modifying the word "empresa" is not considered.

Information about prepositions should be taken into account since they convey important syntactic and se-

mantic information. Let's consider two prepositional expressions: "autorização à empresa" (*permission to the company*) and "autorização da empresa" (*authorization by the company*). According to the Grefenstette's notion of attribute, we should extract the same attribute, namely <*empresa*>, from both expressions. Nevertheless, preposition "a" (*to*) introduces a quite different syntactic dependency than the one introduced by preposition "de" (*by*). Whereas the preposition "a" requires <*empresa*> to be the receiver within the action of giving authorazation, the preposition "de" requires <*empresa*> to be the agent of this action. Therefore, for the purpose of extracting semantic regularities, prepositions should be considered as internal facets of syntactic contexts.

>From the Grefenstette's viewpoint, only one attribute, <*empresa*>, could be extracted from the NN-PREP expression "autorização à empresa". Whereas the modifier word (i.e., the noun after the preposition) is considered as a potential attribute, the modified word (i.e., the noun before the preposition) cannot become an attribute. The tests introduced in section 3.1 will show that the Grefenstette's notion of attribute is too restrictive for the purpose of measuring word similarity. Indeed, the less specific the attribute is, the less precise the word classification will be. In this respect, we should use as specific attributes as possible in order to improve word clustering.

### 2.3 More Accurate Attributes for Word Similarity Measurement

In order to take into account the implicit information contained in the dependency relationships, we will introduce a more general and flexible definition of attribute. The results of the computational tests presented in the next section will provide us with empirical evidence about the appropriateness of such a definition.

Attributes are extracted from binary syntactic dependencies. A syntactic dependency may be represented as the following binary predication:

$$r(w1^{\downarrow}, w2^{\uparrow})$$

this binary predication is constituted by the following entities:

- the binary predicate $r$, which can be associated to specific prepositions, subject relations, direct object relations, etc. ;

- the roles of the predicate, "$\downarrow$" and "$\uparrow$", which represent the *modified* and *modifier* roles, respectively;

- the two words holding the binary relation: w1 and w2.

In this binary syntactic dependency, the word indexed by "$\downarrow$" plays the role of *modified*, whereas the word indexed by "$\uparrow$" plays the role of *modifier*. Therefore, w1 is modified by w2 as well as w2 modifies w1. This way, two complementary attributes may be extracted from that syntactic dependency:

$$< \downarrow r, w1 >< \uparrow r, w2 >$$

where $< \downarrow r, w1 >$ is the attribute of w2 and $< \uparrow r, w2 >$ is the attribute of w1. An attribute is defined as the pair constituted by both a specific syntactic function and the word associated to this function. In particular, $\downarrow r$ represents the syntactic function of *modified*, and $\uparrow r$ the *modifier* function. Consider Table 3. The left column contains expressions constituted by two words syntactically related by a particular type of syntactic dependency. The right column contains the attributes extracted from these expressions. For instance, from the expression "autorização à empresa", it was extracted both the attribute <$\uparrow a$, *empresa*>, where "empresa" plays the role of modifier word, and the attribute <$\downarrow a$, *autorização*>, where "autorização" is the modified word. Let´s note that even though "autorização à empresa" is not truly described as a syntactic constituent by the standard syntagmatic grammar, it should be considered as a very informative syntactic context. Furthermore, information about the specific preposition connecting the words is also available. Our notion of attribute is closely related to what Lin calls "feature" (4).

These elaborate attributes provide us with fine-grained syntactic contexts. In the following section, we will compare these informative syntactic contexts to the coarse-grained contexts used by Grefenstette. This will lead us to assume that the elaborate information conveyed by our notion of attribute is able to contribute more accurately to design a suitable strategy for clustering similar words.

## 3 Comparing Syntactic-Based Strategies

Various semantic extraction techniques were applied to the Brazilian Portuguese corpus from NILC (Interinstitutional Center of Computational Linguistics - USP/São Carlos/Brazil), which is constituted by more than 1,400,000 word occurrences. The corpus is analysed by the parser presented in (5). Similarity was computed by measuring the syntactic information shared by 12,359 different nouns on the basis of 32,293 different attributes.

### 3.1 The Weighted Jaccard Similarity Measure

To compare the syntactic contexts of two words, we used as similarity measure a weighted version of the

binary Jaccard measure (7).[2] The binary Jaccard measure, noted BJ, calculates the similarity value between two words, $m$ and $n$, by comparing the attributes they share and do not share:

$$BJ(w_m, w_n) = \frac{|\{w_m \, atts \cap w_n \, atts\}|}{|\{w_m \, atts \cup w_n \, atts\}|}$$

The weighted Jaccard measure considers a global and a local weight for each attribute. The global weight $gw$ takes into account how many different words are associated with a given attribute. It is computed by the following formula:

$$gw(att_j) = 1 - \sum_i \frac{|p_{ij} \log(p_{ij})|}{nrels}$$

where

$$p_{ij} = \frac{freq \, of \, att_j \, with \, w_i}{total \, of \, atts \, for \, w_i}$$

and $nrels$ is the total number of relations extracted from the corpus. The local weight $lw$ is based on the frequency of the attribute with a given word, and it is calculated by:

$$lw(w_i, att_j) = \log(freq \, of \, att_j \, with \, w_i)$$

The whole weight $W$ of an attribute is the multiplication of both the global and the local weights. So, the weighted Jaccard similarity WJ between two words $m$ and $n$ is computed by:

$$WJ(w_m, w_n) = \frac{\sum_j \min(W(w_m, att_j), W(w_n, att_j))}{\sum_j \max(W(w_m, att_j), W(w_n, att_j))}$$

By computing the similarity measure of all word pairs in the corpus, we extracted the list of the most similar words to each word in the corpus. This process was repeated considering different types of syntactic contexts. On the one hand, we tested the relevance of the use of the prepositional information for the attributes' definition. For this purpose, we compared the results obtained from two strategies: "$+prep$–strategy" and "$-prep$–strategy". The former uses attributes containing information about the specific prepositions, while the latter does not use that information. On the other hand, we tested the adequacy of the "↓–attributes" extracted from prepositional dependencies between two noun phrases. For this purpose, we also compared two different methods: "↑↓–strategy" and "↑–strategy". The former contains both types of attributes, while the later uses only ↑–attributes.

---

[2]We implemented various statistical measures: coefficient of Jaccard, a different version of the weighted Jaccard, and the particular coefficient of Lin. They did not improve, though, the results obtained from the weighted Jaccard measure described in this section.

## 3.2   $+prep$–**strategy** *versus* $-prep$–**strategy**

We tested first the contribution of the specific prepositions to measure word similarity. The manual evaluation is based on a list of 50 randomly chosen words. The results obtained from both strategies, $+prep$–strategy and $-prep$–strategy, showed that there are no significative differences for the words sharing a great number of attributes (namely, more than 100 attributes). That is, the results are not significatively different for words frequently appearing in the corpus. Nevertheless, when the words sharing less than 100 attributes (in fact, the most abundant in the corpus) was compared, we observed that the lists obtained from the $+prep$–strategy are semanticcally more homogeneous than the lists obtained from the $-prep$–strategy. The consideration of the prepositions in the attributes become these ones less common, incresing their weights. On the other side, when we don't consider prepositions in the attributes, the number of attributes shared by words increases, because different attributes (if considering prepositions) can be seen as the same attribute. But, when the number of shared attributes is very high, these modifications became unrelevant. Table 4 shows some of the lists yielded by these strategies for less frequently appearing words.[3]

These results deserve special comments. Let's take the lists obtained from the word "verdade" (*thuth*). In the $+prep$–strategy, the attribute <↑*em, campo*> (<↑*in, field*>) is shared by "verdade" and "impressão" (*impression*). As its global weight is higher, this attribute make the two words more similar. On the contrary, in the $-prep$–strategy, the no–prepositional attribute <↑, *campo*> has a lower weight, which makes the attribute less significant when computing the similarity between "verdade" and "impressão".

Let's consider another example: the lists obtained from the word "punição" (*punishment*). In the $+prep$–strategy, the attribute <↓*de, efeito*> (<↓*of, effect*>) is shared by the words "punição" and "pena" (*penalty*). Given that its global weight is very high, it contributes to make these words semantically close. On the contrary, in the $-prep$–strategy, the no–prepositional attribute <↓, *efeito*> has a lower weight and, consequently, it cannot be considered as a significant clue when comparing the similarity between the words, so "pena" don't appear in list of most similar words of "punição".

Therefore, it can be assumed that the information about specific prepositions is relevant to characterise and identify the significant syntactic contexts used for the measurement of word similarity.

---

[3]We do not use a systematic evaluation methodology based on machine readable dictionaries or electronic thesaurus, because this sort of lexical resources for Portuguese are not available yet.

### 3.3 ↑↓–strategy *versus* ↑–strategy

We also tested the contribution of the ↓–attributes (extracted from noun phrases) to yield lists of similar words. The lists obtained from ↑↓–strategy are significantly more accurate than those obtained from the ↑–strategy, even for the frequently appearing words. Table 5 illustrates some of the lists extracted from both strategies.

On the basis of the results illustrated above, it can be assumed that the use of ↓–attributes to yield lists of similar words is extremely significant. Indeed, this type of attributes somehow provides information concerning the semantic word class. Consider the ↓–attributes <↓*de, equipe*> (<↓*of, team*>), <↓*de, jogador*> (<↓*of, player*>) and <↓*de, partida*> (<↓*of, match*>), shared by the words "tênis" (*tennis*) and "vôlei" (*volleyball*). As those attributes require nouns denoting the same class, namely *sports*, they can be conceived as syntactic patterns imposing the same selectional restrictions to nouns. Consequently, the nouns appearing with those specific ↓–attributes should belong to the class of sports.

In the Appendix, we compare the lists extracted by using the fine-grained techniques (i.e., both ↑↓–strategy and +*prep*–strategy) to the lists extracted by using the coarse-grained methods: ↑–strategy and −*prep*–strategy. The words constituting the lists obtained from the more informative strategies are semantically more homogeneous than those obtained from the less informative ones.

## 4   Final Remarks

According to the results of the tests described above, the strategies based on rich syntactic contexts are more accurate for the measurement of similar words. Experimental tests demonstrated that similarity measures relied on the fine-grained syntactic contexts we have defined in this paper perform better than those based on poorly defined contexts. This is not surprising since the specific syntactic data that we used to refine syntactic attributes allows us to identify more informative syntactic-semantic dependencies between words. Special attention was paid to the very informative syntactic pattern "N-PREP", even if it was not considered as a syntactic constituent in standard syntagmatic grammar.

Nevertheless, all the syntactic-based approaches (fine-grained and coarse-grained syntactic strategies) are confronted with two sorts of linguistic phenomena: both polysemic words and odd attachments of syntactic dependencies between words.

The lists of words recognised as being similar to a particular word can be semantically heterogeneous because of the lexical polysemy of the compared word. For example, the word "segundo" (*second*) appears to be similar to words describing time such as "minuto" (*minute*), "instante" (*instant*), as well as to words describing sequence such as "quarto" (*fourth*). It shares with the former group attributes requiring a quantity of time (e.g., <↑*de, atraso*> (<↑*of, delay*>), <↑*de, desconto*> (<↑*of, discount*>)), and shares with the later group attributes that could require ordinal modifiers: <↓*adj, mundial*> (<↓*adj, world competition*>). Various clustering proposals have been made so as to group the similar words together along sense axes (2; 7). However, the implementation of a efficient clustering method used for the groupping of semantically homogenous words remains to be future work.

Finally, attachment errors inherited from the parser should be taken into account. The number of these errors increases when the language analysed is not as syntactically rigid as English. The rate of incorrect attachments is probably related to the non predictable constraints on the Portuguese syntactic order. To palliate the noisy information inherited from the parser limitations for Portuguese, we are impelled to find more robust word similarity strategies than those used for English texts. So, the fine-grained strategies defined in this paper could be perceived as an attempt to partially palliate the poor results obtained by parsing Portuguese text corpora.

## References

Carolyn J. Crouch and Bokyung Yang. 1992. Experiments in automatic statistical thesaurus construction. *5th Annual International Conference on Research and Development in Information Retrieval*. Copenhagen, pp. 77-88.

David Faure and Claire N'edellec. 1998. ASIUM: Learning subcategorization frames and restrictions of selection. *10th European Conference on Machine Learning, ECML98*. Workshop on Text Mining.

David Yarowsky. 1992. Word-sense disambiguation using statistical models of Roget's categories trained on large corpora. *Proceedings of 14th International Conference on Computational Linguistics, COLING-92*. Nantes, pp. 454-460.

Dekang Lin. 1998. Automatic Retrieval and Clustering of Similar Words. *Proceedings of the COLING-ACL'98*. Montreal.

Eckhard Bick. 2000. *Portuguese Syntax*. VISL Project. p. 114.

Francesc Ribas Framis. 1995. On Learning More Appropriate Selectional Restrictions. *Proceedings of the 7th Conference of the European Chapter of the Association for Computational Linguistics*. Dublin.

Gregory Grefenstette. 1994. *Explorations in Automatic Thesaurus Discovery*. Kluwer Academic Publishers, USA, p. 305.

Gregory Grefenstette 1995. Evaluation Techniques for Automatic Semantic Extraction: Comparing Sy ntatic and Window Based Approaches. *Corpus processing for Lexical Aquisition*. Ed. Branimir Boguraev and James Pustejovsky, The MIT Press, pp. 205-216.

Philip Resnik. 1999. Semantic Similarity in Taxonomy: An Information-Based Measure and its Application to Problems of Ambiguity in Natural Language. *Journal of Artificial Intelligence Research*. Vol. 11, pp. 95-130.

Young Park, Young Han and Key-Sun Choi. 1995. Automatic thesaurus construction using bayesian networks. *International Conference on Information and Knowledge Management*. Baltimore, pp. 212-217.

## Appendix

| Expressions | Binary Dependencies |
|---|---|
| possible causes | *<ADJ, cause, possible>* |
| the cause of neonatal jaundice | *<NNPREP, cause, jaundice>* |
| no cause could be determined | *<DOBJ, determine, cause>* |
| death cause | *<NN, cause, death>* |

Tabela 1: Exemples of syntactic dependencies

| Binary Dependencies | Attributes of `cause` |
|---|---|
| *<ADJ, cause, possible>* | *<possible>* |
| *<NNPREP, cause, jaundice>* | *<jaundice>* |
| *<DOBJ, determine, cause>* | *<DOBJ, determine>* |
| *<NN, cause, death>* | *<death>* |

Tabela 2: Attributes of `cause`

| Binary Expressions | Attributes |
|---|---|
| autorização à empresa (*permission to the company*) | *<↑a, empresa>*, *<↓a, autorização>* |
| nomeação do presidente (*appointment of the president*) | *<↑de, presidente>*, *<↓de, nomeação>* |
| nomeou o presidente (*appointed the president*) | *<↑dobj, presidente>*, *<↓dobj* nomear> |
| discutiu sobre a nomeação (*disscussed about the appointment*) | *<↑sobre, nomeação>*, *<↓sobre,* discutir> |

Tabela 3: Elaborate attributes

| Word | Cluster of similar words | |
|---|---|---|
| | $+prep$–strategy | $-prep$–strategy |
| verdade (*truth*) | desafio, impressão, notícia, dado, responsabilidade (*challenge, impression, news, datum, responsability*) | desafio, culpa, Bélgica, impressão, notícia (*challenge, guilt, Belgium, impression, news*) |
| velocidade (*speed*) | ritmo, vantagem, nível, força, pressão (*rhythm, advantage, level, power, pression*) | vantagem, ritmo, nível, força, diferença (*advantage, rhythm, level, power, difference*) |
| valor (*value*) | preço, salário, proposta, índice, ritmo (*price, salary, proposal, index, rhythm*) | salário, proposta, preço, índice, coisa (*salary, proposal, price, index, thing*) |
| turno (*turn*) | set, rodada, fase, Copa do Mundo, returno (*set, round, phase, World Cup, return*) | Copa do Mundo, rodada, fase, set, Campeonato Brasileiro (*World Cup, round, phase, set, Brazilian Championship*) |
| tragédia (*tragedy*) | zebra, polêmica, milagre, terremoto, ferimento (*strange event, polemics, miracle, earthquake, injury*) | zebra, kart, desperdício, superfície, ferimento (*strange event, kart, wastefulness, surface, injury*) |
| tática (*tactics*) | sistema, fundamento, esquema, pressão, rendimento (*system, foundations, scheme, pression, income*) | fundamento, trajeto, desentendimento, sistema, circunstância (*foundations, way, misunderstanding, system, circunstance*) |
| talento (*talent*) | potencial, facilidade, craque, criatividade, revelação (*potential, facility, excellent player, creativity, revelation*) | reforço, potencial, movimentação, liberdade, facilidade (*reinforcement, potential, movement, freedom, facility*) |
| região (*region*) | centro, interior, litoral, oeste, cidade (*center, countryside, seaside, west, city*) | centro, cidade, mestre, estado, oeste (*center, city, master, state, west*) |
| punição (*punishment*) | suspensão, destaque, pena, atenção, briga (*suspension, prominence, penalty, attention, strife*) | suspensão, destaque, obrigação, mudança, movimentação (*suspension, prominence, obligation, change, movement*) |
| plano (*plan*) | esquema, programa, trabalho, sistema, objetivo (*scheme, program, work, system, objective*) | esquema, amistoso, trabalho, dinheiro, programa (*scheme, amicable game, work, money, program*) |

Tabela 4: Similarity lists of less frequently appearing words ($< 100$ attributes) produced by contexts with and without prepositional information

| Word | Cluster of similar words | |
|---|---|---|
| | ↑↓–strategy | ↑–strategy |
| zaga | defesa, zagueiro, goleiro, ataque, meio-campo | diretoria, meio-campo, fracasso, defesa, área |
| (*defence position*) | (*defence*, *defence player*, *goal-keeper*, *attack*, *middle-field*) | (*management*, *middle-field*, *failure*, *defence*, *area*) |
| violência | briga, pressão, confusão, festa, segurança | briga, conselho, substituição, zagueiro, cobrança |
| (*violence*) | (*strife*, *pression*, *confusion*, *party*, *safety*) | (*strife*, *counsel*, *substitution*, *defence player*, *exaction*) |
| vencedor | campeão, desafio, Japão, equipe, Grécia | Japão, Cuba, campeão, desafio, Grécia |
| (*winner*) | (*champion*, *challenge*, *Japan*, *team*, *Greece*) | (*Japan*, *Cuba*, *champion*, *challenge*, *Greece*) |
| TV | televisão, imprensa, TVA, revista, jornal | imprensa, TVA, revista, jornal, mapa |
| (*TV*) | (*television*, *press*, *TVA*, *magazine*, *newspaper*) | (*press*, *TVA*, *magazine*, *newspaper*, *map*) |
| tênis | basquete, vôlei, surfe, liga, boxe | basquete, surfe, recordista, promessa, mapa |
| (*tennis*) | (*basketball*, *voleyball*, *surf*, *league*, *boxe*) | (*basketball*, *surf*, *record holder*, *promise*, *map*) |
| surpresa | novidade, preocupação, destaque, revelação, atração | preocupação, destaque, novidade, responsabilidade, revelação |
| (*surprise*) | (*novelty*, *preoccupation*, *prominence*, *revelation*, *attraction*) | (*preoccupation*, *prominence*, *novelty*, *responsability*, *revelation*) |
| surfe | tênis, amador, WCT, vôlei, bicampeão | cinema, recordista, amador, tênis, bicampeonato |
| (*surf*) | (*tennis*, *amateur*, *WCT*, *volleyball*, *champion*) | (*cinema*, *record holder*, *amateur*, *tennis*, *championship*) |
| sessão | reunião, fila, evento, divisão, etapa | divisão, fila, turno, semana, evento |
| (*session*) | (*meeting*, *queue*, *event*, *division*, *degree*) | (*division*, *queue*, *turn*, *week*, *event*) |
| regulamento | regra, fórmula, lei, formação, tabela | fórmula, artilheiro, regra, formação, lei |
| (*regulation*) | (*rule*, *formula*, *law*, *formation*, *table*) | (*formula*, *artilleryman*, *rule*, *formation*, *law*) |
| questão | motivo, coisa, erro, destaque, diferença | segurança, destaque, motivo, erro, trabalho |
| (*question*) | (*reason*, *thing*, *error*, *prominence*, *difference*) | (*safety*, *prominence*, *reason*, *error*, *work*) |
| proposta | alternativa, intenção, oferta, convite, notícia | salário, alternativa, intenção, contrato, valor |
| (*proposal*) | (*alternative*, *intention*, *offer*, *invitation*, *news*) | (*salary*, *alternative*, *intention*, *contract*, *value*) |

Tabela 5: Similarity lists produced by contexts with and without ↓–attributes

| Word | Cluster of similar words | |
|---|---|---|
| | +*prep* –strategy and ↑↓–strategy | −*prep*–strategy and ↑–strategy |
| zaga | defesa, zagueiro, goleiro, ataque, meio-campo | diretoria, meio-campo, fracasso, defesa, zagueiro |
| (*defence position*) | (*defence, defence player, goal-keeper, attack, middle-field*) | (*management, middle-field, failure, defence, defence player*) |
| velocidade | ritmo, vantagem, nível, força, pressão | nível, vantagem, diferença, força, ritmo |
| (*speed*) | (*rhythm, advantage, level, power, pression*) | (*level, advantage, difference, power, rhythm*) |
| vantagem | diferença, espaço, oportunidade, média, chance | espaço, velocidade, ritmo, chance, média |
| (*advantage*) | (*difference, space, oportunity, average, chance*) | (*space, speed, rhythm, chance, average*) |
| uniforme | camiseta, camisa, ataque, calção, verão | ânimo, camiseta, despedida, rivalidade, camisa |
| (*uniform*) | (*T-shirt, shirt, attack, ???, summer*) | (*courage, T-shirt, leave-taking, rivalry, shirt*) |
| TV | televisão, imprensa, TVA, revista, jornal | imprensa, TVA, revista, jornal, Globosat |
| (*TV*) | (*television, press, TVA, magazine, newspaper*) | (*press, TVA, magazine, newspaper, Globosat*) |
| tranquilidade | estrutura, virtude, confiança, facilidade, equilíbrio | potencial, estrutura, sorte, personalidade, atenção |
| (*calm*) | (*structure, virtue, confidence, facility, balance*) | (*potential, structure, lucky, personality, attention*) |
| tiro | chute, forte, finalização, pontapé, passe | pancada, aula, fama, forte, passe |
| (*shot*) | (*shot, powerful, finalization, kick, pass*) | (*stroke, class, fame, powerful, pass*) |
| término | abertura, adiamento, fim, rescisão, começo | realização, rescisão, abertura, returno, fim |
| (*finish*) | (*opening, postponement, end, rescission, beginning*) | (*realization, rescission, opening, return, end*) |
| talento | potencial, ídolo, facilidade, craque, fama | ídolo, movimentação, reforço, proteção, liberdade |
| (*talent*) | (*potential, idol, facility, excellent player, fame*) | (*idol, movement, reinforcement, protection, freedom*) |
| sessão | reunião, fila, evento, divisão, etapa | semestre, divisão, fila, turno, encontro |
| (*session*) | (*meeting, queue, event, division, degree*) | (*semester, division, queue, turn, appointment*) |
| regulamento | regra, fórmula, lei, formação, tabela | fórmula, artilheiro, regra, formação, lei |
| (*regulation*) | (*rule, formula, law, formation, table*) | (*formula, artilleryman, rule, formation, law*) |
| recuperação | preparação, preparo, rendimento, tratamento, reação | rendimento, preparo, desgaste, tratamento, estrutura |
| (*recuperation*) | (*preparation, preparing, yield, treatment, reaction*) | (*yield, preparing, wear, treatment, structure*) |
| questão | motivo, coisa, erro, destaque, diferença | destaque, motivo, coisa, trabalho, segurança |
| (*question*) | (*reason, thing, error, prominence, difference*) | (*prominence, reason, thing, work, safety*) |
| proposta | alternativa, intenção, oferta, convite, notícia | contrato, salário, novidade, intenção, valor |
| (*proposal*) | (*alternative, intention, offer, invitation, news*) | (*contract, salary, novelty, intention, value*) |
| programa | plano, corrida, evento, material, formação | corrida, livro, plano, material, esquema |
| (*program*) | (*plan, race, event, material, formation*) | (*race, book, plan, material, scheme*) |
| ponto | gol, vitória, título, chance, resultado | vitória, título, chance, gol, lugar |
| (*point*) | (*goal, victory, title, chance, result*) | (*victory, title, chance, goal, place*) |
| pai | mãe, família, irmão, filho, criança | Senna, mãe, criança, família, pessoa |
| (*father*) | (*mother, family, brother, son, child*) | (*Senna, mother, child, family, person*) |
| luta | confronto, clássico, briga, corrida, competição | corrida, clássico, confronto, lance, conquista |
| (*fight*) | (*confrontation, classic, strife, race, competition*) | (*race, classic, confrontation, throw, conquest*) |
| estratégia | mentalidade, virtude, ponto fraco, tática, visão | virtude, destino, mentalidade, endereço, sistema |
| (*strategy*) | (*mentality, virtue, weak point, tatics, vision*) | (*virtue, destine, mentality, address, system*) |
| esporte | futebol, coisa, trabalho, vôlei, mudança | vôlei, passatempo, criança, futebol, copa |
| (*sport*) | (*soccer, thing, work, volleyball, change*) | (*volleyball, pastime, child, soccer, cup*) |
| clube | equipe, técnico, seleção, futebol, piloto | piloto, atacante, equipe, técnico, futebol |
| (*club*) | (*team, trainer, selection, soccer, pilot*) | (*pilot, attack player, team, trainer, soccer*) |
| chance | oportunidade, possibilidade, condição, ponto, gol | partida, ponto, gol, vitória, lugar |
| (*chance* ) | (*oportunity, possibility, condition, point, goal*) | (*match, point, goal, victory, place*) |

Tabela 6: Similarity lists produced by two sorts of contexts: contexts with +prep and ↑↓–attributes and contexts with −prep and ↑–attributes (like Grefenstette's strategy).