

# Fully Unsupervised Graph-Based Discovery of General-Specific Noun Relationships from Web Corpora Frequency Counts

Gaël Dias  
HULTIG  
University of  
Beira Interior  
ddg@di.ubi.pt

Raycho Mukelov  
HULTIG  
University of  
Beira Interior  
raicho@hultig.di.ubi.pt

Guillaume Cleuziou  
LIFO  
University of  
Orléans  
cleuziou@univ-orleans.pt

## Abstract.

In this paper, we propose a new methodology based on directed graphs and the TextRank algorithm to automatically induce general-specific noun relations from web corpora frequency counts. Different asymmetric association measures are implemented to build the graphs upon which the TextRank algorithm is applied and produces an ordered list of nouns from the most general to the most specific. Experiments are conducted based on the WordNet noun hierarchy and assess 65.69% of correct word ordering.

## 1 Introduction

Taxonomies are crucial for any knowledge-based system. They are in fact important because they allow to structure information, thus fostering their search and reuse. However, it is well known that any knowledge-based system suffers from the so-called knowledge acquisition bottleneck, i.e. the difficulty to actually model the domain in question. As stated in (Caraballo, 1999), WordNet has been an important lexical knowledge base, but it is insufficient for domain specific texts. So, many attempts have been made to automatically produce taxonomies (Grefenstette, 1994), but (Caraballo, 1999) is certainly the first work which proposes a complete overview of the problem by (1) automatically building a hierarchical structure of nouns based on bottom-up clustering methods and (2) labeling the internal nodes of the resulting tree with hypernyms from the nouns clustered underneath by using patterns such as “B is a kind of A”.

In this paper, we are interested in dealing with the second problem of the construction of an organized lexical resource i.e. discovering general-specific noun relationships, so that correct nouns are chosen to label internal nodes of any hierarchical knowledge base, such as the one proposed in (Dias et al., 2006). Most of the works proposed so far have (1) used predefined patterns or (2) automatically learned these patterns to identify hypernym/hyponym relationships. From the first paradigm, (Hearst, 1992) first identifies a set of lexico-syntactic patterns that are easily recognizable i.e. occur frequently and across text genre boundaries. These can be called seed patterns. Based on these seeds, she proposes a bootstrapping algorithm to semi-automatically acquire new more specific patterns. Similarly, (Caraballo, 1999) uses predefined patterns such as “X is a kind of Y” or “X, Y, and other Zs” to identify hypernym/hyponym relationships. This approach to information extraction is based on a technique called *selective concept extraction* as defined by (Riloff, 1993). Selective concept extraction is a form of text skimming that selectively processes relevant text while effectively ignoring surrounding text that is thought to be irrelevant to the domain.

A more challenging task is to automatically learn the relevant patterns for the hypernym/hyponym relationships. In the context of pattern extraction, there exist many approaches as summarized in (Stevenson and Greenwood, 2006). The most well-known work in this area is certainly the one proposed by (Snow et al., 2005) who use machine learning techniques to automatically replace hand-built knowledge. By using dependency path features extracted from parse trees, they introduce a general-purpose formalization and generalization of these patterns. Given a training set of text containing known hypernym pairs, their algorithm automatically extracts useful dependency paths and applies them to new corpora to identify novel pairs. (Sang and Hof-

mann, 2007) use a similar way as (Snow et al., 2006) to derive extraction patterns for hypernym/hyponym relationships by using web search engine counts from pairs of words encountered in WordNet. However, the most interesting work is certainly proposed by (Bollegala et al., 2007) who extract patterns in two steps. First, they find lexical relationships between synonym pairs based on snippets counts and apply wildcards to generalize the acquired knowledge. Then, they apply a SVM classifier to determine whether a new pair shows a relation of synonymy or not, based on a feature vector of lexical relationships. This technique could be applied to hypernym/hyponym relationships although the authors do not mention it.

On the one hand, links between words that result from manual or semi-automatic acquisition of relevant predicative or discursive patterns (Hearst, 1992; Carballo, 1999) are fine and accurate, but the acquisition of these patterns is a tedious task that requires substantial manual work. On the other hand, works done by (Snow et al., 2005; Snow et al., 2006; Sang and Hofmann, 2007; Bollegala et al., 2007) have proposed methodologies to automatically acquire these patterns mostly based on supervised learning to leverage manual work. However, training sets still need to be built.

Unlike other approaches, we propose an unsupervised methodology which aims at discovering general-specific noun relationships which can be assimilated to hypernym/hyponym relationships detection<sup>2</sup>. The advantages of this approach are clear as it can be applied to any language or any domain without any previous knowledge, based on a simple assumption: specific words tend to attract general words with more strength than the opposite. As (Michelbacher et al., 2007) state: “there is a tendency for a strong forward association from a specific term like *adenocarcinoma* to the more general term *cancer*, whereas the association from *cancer* to *adenocarcinoma* is weak”. Based on this assumption, we propose a methodology based on directed graphs and the TextRank algorithm (Mihalcea and Tarau, 2004) to automatically induce general-specific noun relationships from web corpora frequency counts. Indeed, asymmetry in Natural Language Processing can be seen as a possible reason for

the degree of generality of terms (Michelbacher et al., 2007). So, different asymmetric association measures are implemented to build the graphs upon which the TextRank algorithm is applied and produces an ordered list of nouns, from the most general to the most specific. Experiments have been conducted based on the WordNet noun hierarchy and assessed that 65% of the words are ordered correctly.

## 2 Asymmetric Association Measures

In (Michelbacher et al., 2007), the authors clearly point at the importance of asymmetry in Natural Language Processing. In particular, we deeply believe that asymmetry is a key factor for discovering the degree of generality of terms. It is cognitively sensible to state that when someone hears about *mango*, he may induce the properties of a *fruit*. But, when hearing *fruit*, more common fruits will be likely to come into mind such as *apple* or *banana*. In this case, there exists an oriented association between *fruit* and *mango* (*mango* → *fruit*) which indicates that *mango* attracts more *fruit* than *fruit* attracts *mango*. As a consequence, *fruit* is more likely to be a more general term than *mango*.

Based on this assumption, asymmetric association measures are necessary to induce these associations. (Pecina and Schlesinger, 2006) and (Tan et al., 2004) propose exhaustive lists of association measures from which we present the asymmetric ones that will be used to measure the degree of attractiveness between two nouns,  $x$  and  $y$ , where  $f(.,.)$ ,  $P(.,.)$ ,  $P(.,.)$  and  $N$  are respectively the frequency function, the marginal probability function, the joint probability function and the total of digrams.

$$\text{Braun - Blanquet} = \frac{f(x,y)}{\max(f(x,y)+f(x,\bar{y}), f(x,y)+f(\bar{x},y))} \quad (1)$$

$$\text{J measure} = \max \left[ \begin{array}{l} P(x,y) \log \frac{P(y|x)}{P(y)} + P(x,\bar{y}) \log \frac{P(\bar{y}|\bar{x})}{P(\bar{y})} \\ P(x,y) \log \frac{P(x|y)}{P(x)} + P(\bar{x},y) \log \frac{P(x|y)}{P(x)} \end{array} \right] \quad (2)$$

$$\text{Confidence} = \max[P(x|y), P(y|x)] \quad (3)$$

$$\text{Laplace} = \max \left[ \frac{N \cdot P(x,y) + 1}{N \cdot P(x) + 2}, \frac{N \cdot P(x,y) + 1}{N \cdot P(y) + 2} \right] \quad (4)$$

$$\text{Conviction} = \max \left[ \frac{P(x) \cdot P(\bar{y})}{P(x,\bar{y})}, \frac{P(\bar{x}) \cdot P(y)}{P(\bar{x},y)} \right] \quad (5)$$

<sup>2</sup> We must admit that other kinds of relationships may be covered. For that reason, we will speak about general-specific relationships instead of hypernym/hyponym relationships.

$$\text{Certainty Factor} = \max \left[ \frac{P(y|x) - P(y)}{1 - P(y)}, \frac{P(x|y) - P(x)}{1 - P(x)} \right] \quad (6)$$

$$\text{Added Value} = \max [P(y|x) - P(y), P(x|y) - P(x)] \quad (7)$$

All seven definitions show their asymmetry by evaluating the maximum value between two hypotheses i.e. by evaluating the attraction of  $x$  upon  $y$  but also the attraction of  $y$  upon  $x$ . As a consequence, the maximum value will decide the direction of the general-specific association i.e. ( $x \rightarrow y$ ) or ( $y \rightarrow x$ ).

### 3 TextRank Algorithm

Graph-based ranking algorithms are essentially a way of deciding the importance of a vertex within a graph, based on global information recursively drawn from the entire graph. The basic idea implemented by a graph-based ranking model is that of voting or recommendation. When one vertex links to another one, it is basically casting a vote for that other vertex. The higher the number of votes that are cast for a vertex, the higher the importance of the vertex. Moreover, the importance of the vertex casting the vote determines how important the vote itself is, and this information is also taken into account by the ranking model. Hence, the score associated with a vertex is determined based on the votes that are cast for it, and the score of the vertices casting these votes.

Our intuition of using graph-based ranking algorithms is that more general words will be more likely to have incoming associations as they will be associated to many specific words. On the opposite, specific words will have few incoming associations as they will not attract general words (see Figure 1). As a consequence, the voting paradigm of graph-based ranking algorithms should give more strength to general words than specific ones, i.e. a higher voting score.

For that purpose, we first need to build a directed graph. Informally, if  $x$  attracts more  $y$  than  $y$  attracts  $x$ , we will draw an edge between  $x$  and  $y$  as follows ( $x \rightarrow y$ ) as we want to give more credits to general words. Formally, we can define a directed graph  $G = (V, E)$  with the set of vertices  $V$  (in our case, a set of words) and a set of edges  $E$  where  $E$  is a subset of  $V \times V$  (in our case, defined by the asymmetric association measure value between two words). In Figure 1, we show the directed graph obtained by using the set of words  $V = \{isometry, rate\ of\ growth, growth\ rate, rate\}$  randomly extracted from WordNet where *rate of*

*growth* and *growth rate* are synonyms, *isometry* an hyponym of the previous set and *rate* an hypernym of the same set. The weights associated to the edges have been evaluated by the confidence association measure (Equation 3) based on web search engine counts<sup>3</sup>.

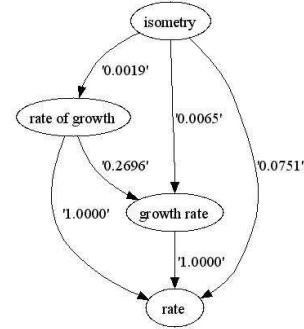


Fig. 1. Directed Graph based on synset #13153496 (*rate of growth, growth rate*) and its direct hypernym (*rate*) and hyponym (*isometry*).

Figure 1 clearly shows our assumption of generality of terms as the hypernym *rate* only has incoming edges whereas the hyponym *isometry* only has outgoing edges. As a consequence, by applying a graph-based ranking algorithm, we aim at producing an ordered list of words from the most general (with the highest value) to the most specific (with the lowest value). For that purpose, we present the TextRank algorithm proposed by (Mihalcea and Tarau, 2004) both for unweighted and weighted directed graphs.

#### 3.1 Unweighted Directed Graph

For a given vertex  $V_i$  let  $In(V_i)$  be the set of vertices that point to it, and let  $Out(V_i)$  be the set of vertices that vertex  $V_i$  points to. The score of a vertex  $V_i$  is defined in Equation 8 where  $d$  is a damping factor that can be set between 0 and 1, which has the role of integrating into the model the probability of jumping from a given vertex to another random vertex in the graph<sup>4</sup>.

$$S(V_i) = (1 - d) + d \times \sum_{V_j \in In(V_i)} \frac{1}{|Out(V_j)|} \times S(V_j) \quad (8)$$

#### 3.2 Weighted Directed Graph

In order to take into account the edge weights, a new formula is introduced in Equation 9.

<sup>3</sup> We used counts returned by <http://www.yahoo.com>.

<sup>4</sup>  $d$  is usually set to 0.85.

$$WS(V_i) = (1 - d) + d \times \sum_{V_j \in In(V_i)} \frac{w_{ji}}{\sum_{V_k \in Out(V_j)} w_{jk}} \times WS(V_j) \quad (9)$$

After running the algorithm in both cases, a score is associated to each vertex, which represents the “importance” of the vertex within the graph. Notice that the final values obtained after TextRank runs to completion are not affected by the choice of the initial values randomly assigned to the vertices. Only the number of iterations needed for convergence may be different. As a consequence, after running the TextRank algorithm, in both its configurations, the output is an ordered list of words from the most general one to the most specific one. In table 1, we show both the lists with the weighted and unweighted versions of the TextRank based on the directed graph shown in Figure 1.

Unweighted		Weighted		WordNet	
$S(V_i)$	Word	$WS(V_i)$	Word	Categ.	Word
0.50	<i>rate</i>	0.81	<i>rate</i>	Hyper.	<i>rate</i>
0.27	<i>growth rate</i>	0.44	<i>growth rate</i>	Synset	<i>growth rate</i>
0.19	<i>rate of growth</i>	0.26	<i>rate of growth</i>	Synset	<i>rate of growth</i>
0.15	<i>isometry</i>	0.15	<i>isometry</i>	Hypo.	<i>isometry</i>

**Table 1.** TextRank ordered lists.

The results show that asymmetric measures combined with directed graphs and graph-based ranking algorithms such as the TextRank are likely to give a positive answer to our hypothesis about the degree of generality of terms. Moreover, we propose an unsupervised methodology for acquiring general-specific noun relationships. However, it is clear that deep evaluation is needed.

## 4 Experiments and Results

Evaluation is classically a difficult task in Natural Language Processing. In fact, as human evaluation is time-consuming and generally subjective even when strict guidelines are provided, measures to automatically evaluate experiments must be proposed. In this section, we propose three evaluation measures and discuss the respective results.

### 4.1 Constraints

WordNet can be defined as applying a set of constraints to words. Indeed, if word  $w$  is the hypernym of word  $x$ , we may represent this relation by the following constraint  $y \succ x$ , where  $\succ$  is the order operator stating that  $y$  is more general than  $x$ . As a consequence, for each set of three

synsets (the hypernym synset, the seed synset and the hyponym synset), a list of constraints can be established i.e. all words of the hypernym synset must be more general than all the words of the seed synset and the hyponym synset, and all the words of the seed synset must be more general than all the words in the hyponym synset. So, if we take the synsets presented in Table 1, we can define the following set of constraints:  $\{rate \succ growth\ rate, rate \succ rate\ of\ growth, growth\ rate \succ isometry, rate\ of\ growth \succ isometry\}$ .

In order to evaluate our list of words ranked by the level of generality against the WordNet categorization, we just need to measure the proportion of constraints which are respected as shown in Equation (10). We call, *correctness* this measure.

$$correctness = \frac{\# \text{ of common constraint}}{\# \text{ of constraint}} \quad (10)$$

For example, in Table 1, all the constraints are respected for both weighted and unweighted graphs, giving 100% correctness for the ordered lists compared to WordNet categorization.

### 4.2 Clustering

Another way to evaluate the quality of the ordering of words is to apply hard clustering to the words weighted by their level of generality. By evidencing the quality of the mapping between three hard clusters generated automatically and the hypernym synset, the seed synset and the hyponym synset, we are able to measure the quality of our ranking. As a consequence, we propose to (1) perform 3-means clustering over the list of ranked words, (2) classify the clusters by level of generality and (3) measure the precision, recall and f-measure of each cluster sorted by level of generality with the hypernym synset, the seed synset and the hyponym synset.

For the first task, we use the implementation of the k-means algorithm of the NLTK toolkit<sup>5</sup>. In particular, we bootstrap the k-means by choosing the initial means as follows. For the first mean, we choose the weight (the score) of the first word in the TextRank generated list of words. For the second mean, we take the weight of the middle word in the list and for the third mean, the weight of the last word in the list.

For the second task the level of generality of each cluster is evaluated by the average level of

<sup>5</sup> <http://nltk.sourceforge.net/>

generality of words inside the cluster (or said with other words by its mean).

For the third task, the most general cluster and the hypernym synset are compared in terms of precision, recall and f-measure as shown in Equation (11), (12) and (13)<sup>6</sup>. The same process is applied to the second most general cluster and the seed synset, and the third cluster and the hyponym synset.

$$precision = \frac{Cluster \cap Synset}{|Cluster|} \quad (11)$$

$$recall = \frac{Cluster \cap Synset}{|Synset|} \quad (12)$$

$$f - measure = \frac{2 \times recall \times precision}{precision + recall} \quad (13)$$

### 4.3 Rank Coefficient Test

The evaluation can be seen as a rank test between two ordered lists. Indeed, one way to evaluate the results is to compare the list of general-specific relationships encountered by the TextRank algorithm and the original list given by WordNet. However, we face one problem. WordNet does not give an order of generality inside synsets. In order to avoid this problem, we can order words in each synset by their estimated frequency given by WordNet<sup>7</sup> as well as their frequency calculated by web search hits. An example of both ordered lists is given in Table 2 for the synset #6655336 and its immediate hypernyms and hyponyms.

WordNet Estimated Frequency		Web Estimated Frequency	
Category	Word	Category	Word
Hypernym	<i>statement</i>	Hypernym	<i>statement</i>
Synset	<i>answer</i>	Synset	<i>reply</i>
Synset	<i>reply</i>	Synset	<i>response</i>
Synset	<i>response</i>	Synset	<i>answer</i>
Hyponym	<i>rescript</i>	Hyponym	<i>feedback</i>
Hyponym	<i>feedback</i>	Hyponym	<i>rescript</i>

**Table 2.** Estimated Frequency ordered lists for synset #6655336.

For that purpose, we propose to use the Spearman’s rank correlation coefficient (Rho). The Spearman’s Rho is a statistical coefficient that shows how much two random variables are cor-

related. It is defined in Equation (14) where  $d$  is the distance between every pair of words in the list ordered with TextRank and the reference list which is ordered according to WordNet or the Web and  $n$  is the number of pairs of ranked words.

$$\rho = 1 - \frac{6 \times \sum d_i^2}{n(n^2 - 1)} \quad (14)$$

In particular, the Spearman’s rank correlation coefficient is a number between -1 (no correlation at all) and 1 (very strong correlation).

## 4.4 Experiments

In order to evaluate our methodology, we randomly<sup>8</sup> extracted 800 seed synsets for which we retrieved their hypernym and hyponym synsets. For each seed synset, we then built the associated directed weighted and unweighted graphs based on the asymmetric association measures referred to in section 2<sup>9</sup> and ran the TextRank algorithm to produce a general-specific ordered lists of terms.

### 4.4.1 Results by Constraints

In Table 3, we present the results of the correctness for all seven asymmetric measures, both for the unweighted and weighted graphs.

Equation	Type of Graph	Correctness
Braun-Blanquet	Unweighted	65.68%
	Weighted	65.52%
J measure	Unweighted	60.00%
	Weighted	60.34%
Confidence	Unweighted	<b>65.69%</b>
	Weighted	65.40%
Laplace	Unweighted	<b>65.69%</b>
	Weighted	<b>65.69%</b>
Conviction	Unweighted	61.81%
	Weighted	63.39%
Certainty Factor	Unweighted	65.59%
	Weighted	63.76%
Added Value	Unweighted	65.61%
	Weighted	64.90%
Baseline <sup>10</sup>	None	55.68%

**Table 3.** Results for the Evaluation by Constraints.

The best results are obtained by the Confidence and the Laplace measures reaching 65.69% cor-

<sup>6</sup> Where  $Cluster \cap Synset$  means the number of words common to both Synset and Cluster, and  $|Synset|$  and  $|Cluster|$  respectively measure the number of words in the Synset and the Cluster.

<sup>7</sup> We use WordNet 2.1.

<sup>8</sup> We guarantee 98% significance level for an error of 0.05 following the normal distribution.

<sup>9</sup> The probability functions are estimated by the Maximum Likelihood Estimation (MLE).

<sup>10</sup> The baseline is the list of words ordered by web hits frequency (without TextRank).

rectness. However, the Braun-Blanquet, the Certainty Factor and the Added Value give results near the best ones. Only the J measure and the Conviction metric seem to perform worst.

It is also important to note that the difference between unweighted and weighted graphs is marginal which clearly points at the fact that the topology of the graph is more important than its weighting. This is also confirmed by the fact that most of the asymmetric measures perform alike.

#### 4.4.2 Results by Clustering

In Table 4, we present the results of precision, recall and f-measure for both weighted and unweighted graphs for all the seven asymmetric measures. The best precision is obtained for the weighted graph with the Confidence measure evidencing 47.62% and the best recall is also obtained by the Confidence measure also for the weighted graph reaching 47.68%. Once again, the J measure and the Conviction metric perform worst showing worst f-measures. Contrarily, the Confidence measure shows the best performance in terms of f-measure for the weighted graph, i.e. 47.65% while the best result for the unweighted graphs is obtained by the Certainty factor with 46.50%.

These results also show that the weighting of the graph plays an important issue in our methodology. Indeed, most metrics perform better with weighted graphs in terms of f-measure.

Equation	Graph	Precision	Recall	F-measure
Braun-Blanquet	Unweighted	46.61	46.06	46.33
	Weighted	47.60	47.67	47.64
J measure	Unweighted	40.92	40.86	40.89
	Weighted	42.61	43.71	43.15
Confidence	Unweighted	46.54	46.02	46.28
	Weighted	<b>47.62</b>	<b>47.68</b>	<b>47.65</b>
Laplace	Unweighted	<b>46.67</b>	46.11	46.39
	Weighted	46.67	46.11	46.39
Conviction	Unweighted	42.13	41.67	41.90
	Weighted	43.62	43.99	43.80
Certainty Factor	Unweighted	46.49	<b>46.52</b>	<b>46.50</b>
	Weighted	44.84	45.85	45.34
Added Value	Unweighted	46.61	46.59	46.60
	Weighted	47.13	47.27	47.19

**Table 4.** Results for the Evaluation by Clustering.

In Table 5, 6 and 7, we present the same results as in Table 4 but at different levels of analysis i.e. precision, recall and f-measure at hypernym, seed and hyponym levels. Indeed, it is important to understand how the methodology performs at different levels of generality as we verified that

our approach performs better at higher levels of generality.

Equation	Graph	Precision	Recall	F-measure
Braun-Blanquet	Unweighted	59.38	37.38	45.88
	Weighted	58.75	39.35	47.14
J measure	Unweighted	46.49	37.00	41.20
	Weighted	47.19	41.90	44.38
Confidence	Unweighted	59.20	37.30	45.77
	Weighted	58.71	39.22	47.03
Laplace	Unweighted	<b>59.50</b>	37.78	<b>45.96</b>
	Weighted	<b>59.50</b>	37.78	45.96
Conviction	Unweighted	50.07	35.88	41.80
	Weighted	52.72	40.74	45.96
Certainty Factor	Unweighted	55.90	<b>38.29</b>	45.45
	Weighted	51.64	<b>42.93</b>	46.88
Added Value	Unweighted	56.26	37.90	45.29
	Weighted	58.21	40.09	<b>47.48</b>

**Table 5.** Results at the hypernym level.

Equation	Graph	Precision	Recall	F-measure
Braun-Blanquet	Unweighted	43.05	37.86	40.29
	Weighted	<b>46.38</b>	33.14	38.66
J measure	Unweighted	40.82	<b>43.72</b>	42.22
	Weighted	43.98	33.89	38.28
Confidence	Unweighted	43.03	37.67	40.17
	Weighted	46.36	33.02	38.57
Laplace	Unweighted	43.10	37.78	40.27
	Weighted	43.10	37.78	40.27
Conviction	Unweighted	40.36	38.02	39.16
	Weighted	42.60	26.39	32.59
Certainty Factor	Unweighted	<b>44.28</b>	40.87	<b>42.51</b>
	Weighted	44.14	<b>40.70</b>	<b>42.35</b>
Added Value	Unweighted	44.21	40.74	42.40
	Weighted	45.78	32.90	38.29

**Table 6.** Results at the seed level.

Equation	Graph	Precision	Recall	F-measure
Braun-Blanquet	Unweighted	37.39	62.96	46.92
	Weighted	37.68	70.50	49.12
J measure	Unweighted	35.43	41.87	38.38
	Weighted	36.69	55.33	44.12
Confidence	Unweighted	37.38	63.09	46.95
	Weighted	37.79	<b>70.80</b>	<b>49.27</b>
Laplace	Unweighted	37.40	<b>63.11</b>	46.97
	Weighted	37.40	63.11	46.97
Conviction	Unweighted	35.97	50.94	42.16
	Weighted	35.54	64.85	45.92
Certainty Factor	Unweighted	39.28	60.40	47.60
	Weighted	<b>38.74</b>	53.92	45.09
Added Value	Unweighted	<b>39.36</b>	61.15	<b>47.89</b>
	Weighted	37.39	68.81	48.45

**Table 7.** Results at the hyponym level.

Indeed, the precision scores go down from 59.50% at the hypernym level to 39.36% at the hyponym level with 46.38% at the seed level. The same phenomenon is inversely true for the recall with 42.93% at the hypernym level,

43.72% at the seed level and 70.80% at the hyponym level.

This situation can easily be understood as most of the clusters created by the k-means present the same characteristics i.e. the upper level cluster usually has fewer words than the middle level cluster which in turn has fewer words than the last level cluster. As a consequence, the recall is artificially high for the hyponym level. But on the opposite, the precision is high for higher levels of generality which is promising for the automatic construction of hierarchical thesauri. Indeed, our approach can be computed recursively so that each level of analysis is evaluated as if it was at the hypernym level, thus taking advantage of the good performance of our approach at upper levels of generality<sup>11</sup>.

#### 4.4.3 Results by Rank Test

For each produced list, we calculated the Spearman's Rho both with WordNet and Web Estimated Lists for weighted and unweighted graphs. Table 8 presents the average results for the 800 randomly selected synsets.

Equation	Type of Graph	Rho with WNet Est. list	Rho with Web Est. list
Braun-Blanquet	Unweighted	0.38	0.30
	Weighted	<b>0.39</b>	<b>0.39</b>
J measure	Unweighted	0.23	0.19
	Weighted	0.27	0.27
Confidence	Unweighted	0.38	0.30
	Weighted	<b>0.39</b>	<b>0.39</b>
Laplace	Unweighted	0.38	0.30
	Weighted	0.38	0.38
Conviction	Unweighted	0.30	0.22
	Weighted	0.33	0.33
Certainty Factor	Unweighted	0.38	0.29
	Weighted	0.35	0.35
Added Value	Unweighted	0.37	0.29
	Weighted	0.38	0.38
Baseline <sup>12</sup>	None	0.14	0.14

**Table 8.** Results for the Spearman's rank correlation coefficient.

Similarly to what we evidenced in section 4.4.1., the J measure and the Conviction metric are the measures which less seem to map the correct order by evidencing low correlation scores. On the other hand, the Confidence metric still gives the best results equally with the Laplace and Braun-Blanquet metrics.

<sup>11</sup> This will be studied as future work.

<sup>12</sup> The baseline is the list of words ordered by web hits frequency.

It is interesting to note that in the case of the web estimated list, the weighted graphs evidence much better results than the unweighted ones, although they do not show improved results compared to the WordNet list. On the one hand, these results show that our methodology is capable to map to WordNet lists as easily as to Web lists even that it is based on web frequency counts. On the other hand, the fact that weighted graphs perform best, shows that the topology of the graph lacks in accuracy and needs the application of weights to counterpoint this lack.

#### 4.5 Discussion

An important remark needs to be made at this point of our explanation. There is a large ambiguity introduced in the methodology by just looking at web counts. Indeed, when counting the occurrences of a word like *answer*, we count all its occurrences for all its meanings and forms. For example, based on WordNet, the word *answer* can be a verb with ten meanings and a noun with five meanings. Moreover, words are more frequent than others although they are not so general, unconfirming our original hypothesis. Looking at Table 2, *feedback* is a clear example of this statement. As we are not dealing with a single domain within which one can expect to see the "one sense per discourse" paradigm, it is clear that the Rho coefficient would not be as good as expected as it is clearly biased by "incorrect" counts. One direct implication of this comment is the use of web estimated lists to evaluate the methodology.

Also, there has been a great discussion over the last few months in the corpora list<sup>13</sup> whether one should use web counts instead of corpus counts to estimate word frequencies. In our study, we clearly see that web counts show evident problems, like the ones mentioned by (Kilgarriff, 2007). However, they cannot be discarded so easily. In particular, we aim at looking at web counts in web directories that would act as specific domains and would reduce the space for ambiguity. Of course, experiments with well-known corpora will also have to be made to understand better this phenomenon.

#### 5 Conclusions and Future Work

In this paper, we proposed a new methodology based on directed weighted/unweighted graphs and the TextRank algorithm to automatically in-

<sup>13</sup> Finalized by (Kilgarriff, 2007).

duce general-specific noun relationships from web corpora frequency counts. To our knowledge, such an unsupervised experiment has never been attempted so far. In order to evaluate our results, we proposed three different evaluation metrics. The results obtained by using seven asymmetric association measures based on web frequency counts showed promising results reaching levels of (1) constraint coherence of 65.69%, (2) clustering mapping of 59.50% in terms of precision for the hypernym level and 42.72% on average in terms of f-measure and (3) ranking similarity of 0.39 for the Spearman's rank correlation coefficient.

As future work, we intend to take advantage of the good performance of our approach at the hypernym level to propose a recursive process to improve precision results over all levels of generality.

Finally, it is important to notice that the evaluation by clustering evidences more than a simple evaluation of the word order, but shows how this approach is capable to automatically map clusters to WordNet classification.

## References

- Bollegala, D., Matsuo, Y. and Ishizuka, M. 2007. *Measuring Semantic Similarity between Words Using WebSearch Engines*. In Proceedings of International World Wide Web Conference (WWW 2007).
- Caraballo, S.A. 1999. *Automatic Construction of a Hypernym-labeled Noun Hierarchy from Text*. In Proceedings of the Conference of the Association for Computational Linguistics (ACL 1999).
- Dias, G., Santos, C., and Cleuziou, G. 2006. *Automatic Knowledge Representation using a Graph-based Algorithm for Language-Independent Lexical Chaining*. In Proceedings of the Workshop on Information Extraction Beyond the Document associated to the Joint Conference of the International Committee of Computational Linguistics and the Association for Computational Linguistics (COLING/ACL), pages. 36-47.
- Grefenstette, G. 1994. *Explorations in Automatic Thesaurus Discovery*. Kluwer Academic Publishers, USA.
- Hearst, M.H. 1992. *Automatic Acquisition of Hyponyms from Large Text Corpora*. In Proceedings of the Fourteenth International Conference on Computational Linguistics (COLING 1992), pages 539-545.
- Kilgarriff, A. 2007. Googleology is Bad Science. *Computational Linguistics* 33 (1), pages: 147-151.
- Michelbacher, L., Evert, S. and Schütze, H. 2007. *Asymmetric Association Measures*. In Proceedings of the Recent Advances in Natural Language Processing (RANLP 2007).
- Mihalcea, R. and Tarau, P. 2004. *TextRank: Bringing Order into Texts*. In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP 2004), pages 404-411.
- Pecina, P. and Schlesinger, P. 2006. *Combining Association Measures for Collocation Extraction*. In Proceedings of the International Committee of Computational Linguistics and the Association for Computational Linguistics (COLING/ACL 2006).
- Riloff, E. 1993. *Automatically Constructing a Dictionary for Information Extraction Tasks*. In Proceedings of the Eleventh National Conference on Artificial Intelligence (AAAI 1993), pages 811-816.
- Sang, E.J.K. and Hofmann, K. 2007. *Automatic Extraction of Dutch Hypernym-Hyponym Pairs*. In Proceedings of Computational Linguistics in the Netherlands Conference (CLIN 2007).
- Snow, R., Jurafsky, D. and Ng, A. Y. 2005. *Learning Syntactic Patterns for Automatic Hypernym Discovery*. In Proceedings of the International Committee of Computational Linguistics and the Association for Computational Linguistics (COLING/ACL 2006).
- Snow, R., Jurafsky, D. and Ng, A. Y. 2005. *Semantic Taxonomy Induction from Heterogenous Evidence*. In Proceedings of the Neural Information Processing Systems Conference (NIPS 2005).
- Stevenson, M., and Greenwood, M. 2006. *Comparing Information Extraction Pattern Models*. In Proceedings of the Workshop on Information Extraction Beyond the Document associated to the Joint Conference of the International Committee of Computational Linguistics and the Association for Computational Linguistics (COLING/ACL 2006), pages. 29-35.
- Tan, P.-N., Kumar, V. and Srivastava, J. 2004. *Selecting the Right Objective Measure for Association Analysis*. *Information Systems*, 29(4). pages 293-313.