

TextRank: Bringing Order into Texts

Rada Mihalcea and Paul Tarau
Department of Computer Science
University of North Texas
`{rada,tarau}@cs.unt.edu`

Abstract

In this paper, we introduce TextRank – a graph-based ranking model for text processing, and show how this model can be successfully used in natural language applications. In particular, we propose two innovative unsupervised methods for keyword and sentence extraction, and show that the results obtained compare favorably with previously published results on established benchmarks.

1 Introduction

Graph-based ranking algorithms like Kleinberg’s HITS algorithm (Kleinberg, 1999) or Google’s PageRank (Brin and Page, 1998) have been successfully used in citation analysis, social networks, and the analysis of the link-structure of the World Wide Web. Arguably, these algorithms can be singled out as key elements of the paradigm-shift triggered in the field of Web search technology, by providing a Web page ranking mechanism that relies on the collective knowledge of Web architects rather than individual content analysis of Web pages. In short, a graph-based ranking algorithm is a way of deciding on the importance of a vertex within a graph, by taking into account global information recursively computed from the entire graph, rather than relying only on local vertex-specific information.

Applying a similar line of thinking to lexical or semantic graphs extracted from natural language documents, results in a graph-based ranking model that can be applied to a variety of natural language processing applications, where knowledge drawn from an entire text is used in making local ranking/selection decisions. Such text-oriented ranking methods can be applied to tasks ranging from automated extraction of keyphrases, to extractive summarization and word sense disambiguation (Mihalcea et al., 2004).

In this paper, we introduce the TextRank graph-based ranking model for graphs extracted from natural language texts. We investigate and evaluate the application of TextRank to two language processing tasks consisting of unsupervised keyword and sen-

tence extraction, and show that the results obtained with TextRank are competitive with state-of-the-art systems developed in these areas.

2 The TextRank Model

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.

Formally, let $G = (V, E)$ be a directed graph with the set of vertices V and set of edges E , where E is a subset of $V \times V$. For a given vertex V_i , let $In(V_i)$ be the set of vertices that point to it (predecessors), and let $Out(V_i)$ be the set of vertices that vertex V_i points to (successors). The score of a vertex V_i is defined as follows (Brin and Page, 1998):

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

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. In the context of Web surfing, this graph-based ranking algorithm implements the “random surfer model”, where a user clicks on links at random with a probability d , and jumps to a completely new page with probability $1 - d$. The factor d is usually set to 0.85 (Brin and Page, 1998), and this is the value we are also using in our implementation.

document->extract sentences-->extract keywords

1. Identify text units that best define the task at hand, and add them as vertices in the graph.
2. Identify relations that connect such text units, and use these relations to draw edges between vertices in the graph. Edges can be directed or undirected, weighted or unweighted.
3. Iterate the graph-based ranking algorithm until convergence.
4. Sort vertices based on their final score. Use the values attached to each vertex for ranking/selection decisions.

In the following, we investigate and evaluate the application of TextRank to two natural language processing tasks involving ranking of text units: (1) A keyword extraction task, consisting of the selection of keyphrases representative for a given text; and (2) A sentence extraction task, consisting of the identification of the most “important” sentences in a text, which can be used to build extractive summaries.

3 Keyword Extraction frequency not working

The task of a keyword extraction application is to automatically identify in a text a set of terms that best describe the document. Such keywords may constitute useful entries for building an automatic index for a document collection, can be used to classify a text, or may serve as a concise summary for a given document. Moreover, a system for automatic identification of important terms in a text can be used for the problem of terminology extraction, and construction of domain-specific dictionaries.

The simplest possible approach is perhaps to use a frequency criterion to select the “important” keywords in a document. However, this method was generally found to lead to poor results, and consequently other methods were explored. The state-of-the-art in this area is currently represented by supervised learning methods, where a system is trained to recognize keywords in a text, based on lexical and syntactic features. This approach was first suggested in (Turney, 1999), where parametrized heuristic rules are combined with a genetic algorithm into a system for keyphrase extraction - GenEx - that automatically identifies keywords in a document. A different learning algorithm was used in (Frank et al., 1999), where a Naive Bayes learning scheme is applied on the document collection, with improved results observed on the same data set as used in (Turney, 1999). Neither Turney nor Frank report on the recall of their systems, but only on precision: a 29.0% precision is achieved with GenEx (Turney, 1999) for five keyphrases extracted per document, and 18.3% precision achieved with Kea (Frank et al., 1999) for fifteen keyphrases per document.

More recently, (Hulth, 2003) applies a supervised learning system to keyword extraction from ab-

stracts, using a combination of lexical and syntactic features, proved to improve significantly over previously published results. As Hulth suggests, keyword extraction from abstracts is more widely applicable than from full texts, since many documents on the Internet are not available as full-texts, but only as abstracts. In her work, Hulth experiments with the approach proposed in (Turney, 1999), and a new approach that integrates part of speech information into the learning process, and shows that the accuracy of the system is almost doubled by adding linguistic knowledge to the term representation.

In this section, we report on our experiments in keyword extraction using TextRank, and show that the graph-based ranking model outperforms the best published results in this problem. Similar to (Hulth, 2003), we are evaluating our algorithm on keyword extraction from abstracts, mainly for the purpose of allowing for a direct comparison with the results she reports with her keyphrase extraction system. Notice that the size of the text is not a limitation imposed by our system, and similar results are expected with TextRank applied on full-texts.

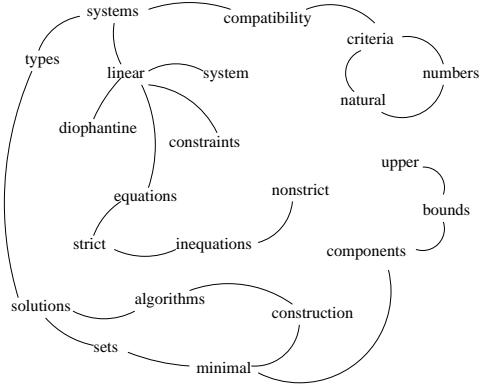
3.1 TextRank for Keyword Extraction

The expected end result for this application is a set of words or phrases that are representative for a given natural language text. The units to be ranked are therefore sequences of one or more lexical units extracted from text, and these represent the vertices that are added to the text graph. Any relation that can be defined between two lexical units is a potentially useful connection (edge) that can be added between two such vertices. We are using a *co-occurrence* relation, controlled by the distance between word occurrences: two vertices are connected if their corresponding lexical units co-occur within a window of maximum N words, where N can be set anywhere from 2 to 10 words. Co-occurrence links express relations between syntactic elements, and similar to the semantic links found useful for the task of word sense disambiguation (Mihalcea et al., 2004), they represent cohesion indicators for a given text.

The vertices added to the graph can be restricted with syntactic filters, which select only lexical units of a certain part of speech. One can for instance consider only nouns and verbs for addition to the graph, and consequently draw potential edges based only on relations that can be established between nouns and verbs. We experimented with various syntactic filters, including: all open class words, nouns and verbs only, etc., with best results observed for nouns and adjectives only, as detailed in section 3.2.

The TextRank keyword extraction algorithm is fully unsupervised, and proceeds as follows. First,

Compatibility of systems of linear constraints over the set of natural numbers. Criteria of compatibility of a system of linear Diophantine equations, strict inequations, and nonstrict inequations are considered. Upper bounds for components of a minimal set of solutions and algorithms of construction of minimal generating sets of solutions for all types of systems are given. These criteria and the corresponding algorithms for constructing a minimal supporting set of solutions can be used in solving all the considered types systems and systems of mixed types.



Keywords assigned by TextRank:
linear constraints; linear diophantine equations; natural numbers; nonstrict inequations; strict inequations; upper bounds
Keywords assigned by human annotators:
linear constraints; linear diophantine equations; minimal generating sets; non-strict inequations; set of natural numbers; strict inequations; upper bounds

Figure 2: Sample graph build for keyphrase extraction from an *Inspec* abstract

the text is tokenized, and annotated with part of speech tags – a preprocessing step required to enable the application of syntactic filters. To avoid excessive growth of the graph size by adding all possible combinations of sequences consisting of more than one lexical unit (ngrams), we **consider only single words as candidates for addition to the graph**, with multi-word keywords being eventually reconstructed in the post-processing phase.

Next, all lexical units that pass the syntactic filter are added to the graph, and an edge is added between those lexical units that co-occur within a window of N words. After the graph is constructed (undirected unweighted graph), the score associated with each vertex is set to an initial value of 1, and the ranking algorithm described in section 2 is run on the graph for several iterations until it converges – usually for 20-30 iterations, at a threshold of 0.0001.

Once a final score is obtained for each vertex in the graph, **vertices are sorted in reversed order of their score, and the top T vertices in the ranking are retained for post-processing**. While T may be set to any fixed value, usually ranging from 5 to 20 keywords (e.g. (Turney, 1999) limits the number of keywords extracted with his GenEx system to five), we are using a more flexible approach, which decides

the number of keywords **based on the size of the text**. For the data used in our experiments, which consists of relatively short abstracts, T is set to a third of the number of vertices in the graph.

During post-processing, all lexical units selected as potential keywords by the TextRank algorithm are marked in the text, and sequences of adjacent keywords are collapsed into a multi-word keyword. For instance, in the text *Matlab code for plotting ambiguity functions*, if both *Matlab* and *code* are selected as potential keywords by TextRank, since they are adjacent, they are collapsed into one single keyword *Matlab code*.

Figure 2 shows a sample graph built for an abstract from our test collection. While the size of the abstracts ranges from 50 to 350 words, with an average size of 120 words, we have deliberately selected a very small abstract for the purpose of illustration. For this example, the lexical units found to have higher “importance” by the TextRank algorithm are (with the TextRank score indicated in parenthesis): numbers (1.46), inequations (1.45), linear (1.29), diophantine (1.28), upper (0.99), bounds (0.99), strict (0.77). Notice that this ranking is different than the one rendered by simple word frequencies. For the same text, a frequency approach provides the following top-ranked lexical units: systems (4), types (3), solutions (3), minimal (3), linear (2), inequations (2), algorithms (2). All other lexical units have a frequency of 1, and therefore cannot be ranked, but only listed.

3.2 Evaluation

The data set used in the experiments is a collection of 500 abstracts from the *Inspec* database, and the corresponding manually assigned keywords. This is the same test data set as used in the keyword extraction experiments reported in (Hulth, 2003). The *Inspec* abstracts are from journal papers from Computer Science and Information Technology. Each abstract comes with two sets of keywords assigned by professional indexers: controlled keywords, restricted to a given thesaurus, and uncontrolled keywords, freely assigned by the indexers. We follow the evaluation approach from (Hulth, 2003), and use the uncontrolled set of keywords.

In her experiments, Hulth is using a total of 2000 abstracts, divided into 1000 for training, 500 for development, and 500 for test². Since our approach is completely unsupervised, no training/development data is required, and we are only using the test docu-

²Many thanks to Anette Hulth for allowing us to run our algorithm on the data set used in her keyword extraction experiments, and for making available the training/test/development data split.

sentence), or longer more explicative summaries, consisting of more than 100 words. We are also investigating combinations of keyphrase and sentence extraction techniques as a method for building short/long summaries.

Finally, another advantage of TextRank over previously proposed methods for building extractive summaries is the fact that it does not require training corpora, which makes it easily adaptable to other languages or domains.

5 Why TextRank Works

Intuitively, TextRank works well because it does not only rely on the local context of a text unit (vertex), but rather it takes into account information recursively drawn from the entire text (graph).

Through the graphs it builds on texts, TextRank identifies connections between various entities in a text, and implements the concept of *recommendation*. A text unit recommends other related text units, and the strength of the recommendation is recursively computed based on the importance of the units making the recommendation. For instance, in the keyphrase extraction application, co-occurring words recommend each other as important, and it is the common context that enables the identification of connections between words in text. In the process of identifying important sentences in a text, a sentence recommends another sentence that addresses similar concepts as being useful for the overall understanding of the text. The sentences that are highly recommended by other sentences in the text are likely to be more informative for the given text, and will be therefore given a higher score.

An analogy can be also drawn with PageRank’s “random surfer model”, where a user surfs the Web by following links from any given Web page. In the context of text modeling, TextRank implements what we refer to as “text surfing”, which relates to the concept of text cohesion (Halliday and Hasan, 1976): from a certain concept C in a text, we are likely to “follow” links to connected concepts – that is, concepts that have a relation with the current concept C (be that a lexical or semantic relation). This also relates to the “knitting” phenomenon (Hobbs, 1974): facts associated with words are shared in different parts of the discourse, and such relationships serve to “knit the discourse together”.

Through its iterative mechanism, TextRank goes beyond simple graph connectivity, and it is able to score text units based also on the “importance” of other text units they link to. The text units selected by TextRank for a given application are the ones most recommended by related text units in the text, with preference given to the recommendations made by

most influential ones, i.e. the ones that are in turn highly recommended by other related units. The underlying hypothesis is that in a cohesive text fragment, related text units tend to form a “Web” of connections that approximates the model humans build about a given context in the process of discourse understanding.

6 Conclusions

In this paper, we introduced TextRank – a graph-based ranking model for text processing, and show how it can be successfully used for natural language applications. In particular, we proposed and evaluated two innovative unsupervised approaches for keyword and sentence extraction, and showed that the accuracy achieved by TextRank in these applications is competitive with that of previously proposed state-of-the-art algorithms. An important aspect of TextRank is that it does not require deep linguistic knowledge, nor domain or language specific annotated corpora, which makes it highly portable to other domains, genres, or languages.

References

- S. Brin and L. Page. 1998. The anatomy of a large-scale hyper-textual Web search engine. *Computer Networks and ISDN Systems*, 30(1–7).
- DUC. 2002. Document understanding conference 2002. <http://www-nplir.nist.gov/projects/duc/>.
- E. Frank, G. W. Paynter, I. H. Witten, C. Gutwin, and C. G. Nevill-Manning. 1999. Domain-specific keyphrase extraction. In *Proceedings of the 16th International Joint Conference on Artificial Intelligence*.
- M. Halliday and R. Hasan. 1976. *Cohesion in English*. Longman.
- P.J. Herings, G. van der Laan, and D. Talman. 2001. Measuring the power of nodes in digraphs. Technical report, Tinbergen Institute.
- J. Hobbs. 1974. A model for natural language semantics. Part I: The model. Technical report, Yale University.
- A. Hulth. 2003. Improved automatic keyword extraction given more linguistic knowledge. In *Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing*, Japan, August.
- J.M. Kleinberg. 1999. Authoritative sources in a hyperlinked environment. *Journal of the ACM*, 46(5):604–632.
- C.Y. Lin and E.H. Hovy. 2003. Automatic evaluation of summaries using n-gram co-occurrence statistics. In *Proceedings of Human Language Technology Conference (HLT-NAACL 2003)*, Edmonton, Canada, May.
- R. Mihalcea, P. Tarau, and E. Figa. 2004. PageRank on semantic networks, with application to word sense disambiguation. In *Proceedings of the 20th International Conference on Computational Linguistics (COLING 2004)*, Geneva, Switzerland.
- R. Mihalcea. 2004. Graph-based ranking algorithms for sentence extraction, applied to text summarization. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL 2004) (companion volume)*, Barcelona, Spain.
- P. Turney. 1999. Learning to extract keyphrases from text. Technical report, National Research Council, Institute for Information Technology.