

# Visualization in Information Retrieval from Hospital Information System

Miroslav Bursa<sup>1</sup>, Lenka Lhotska<sup>1</sup>, Vaclav Chudacek<sup>1</sup>, Jiri Spilka<sup>1</sup>,  
Petr Janku<sup>2</sup>, and Lukas Hruban<sup>1</sup>

<sup>1</sup> Dept. of Cybernetics, Faculty of Electrical Engineering,  
Czech Technical University in Prague, Czech Republic  
`bursam@fel.cvut.cz`

<sup>2</sup> Obstetrics and Gynaecology clinic,  
University Hospital in Brno, Czech Republic

**Abstract.** This paper describes the process of mining information from loosely structured medical textual records with no apriori knowledge. The typical patient record is filled with typographical errors, duplicates, ambiguities, syntax errors and many (nonstandard) abbreviations. In the paper we depict the process of mining a large dataset of ~50,000–120,000 records  $\times$  20 attributes in database tables, originating from the hospital information system (thanks go to the University Hospital in Brno, Czech Republic) recording over 11 years. The proposed technique has an important impact on reduction of the processing time of loosely structured textual records for experts.

Note that this project is an ongoing process (and research) and new data are irregularly received from the medical facility, justifying the need for robust and fool-proof algorithms.

**Keywords:** Swarm Intelligence, Ant Colony, Textual Data Mining, Medical Record Processing, Hospital Information System.

## 1 Introduction

### 1.1 Nature Inspired Methods

Nature inspired metaheuristics play an important role in the domain of artificial intelligence, offering fast and robust solutions in many fields (graph algorithms, feature selection, optimization, clustering, feature selection, etc). Stochastic nature inspired metaheuristics have interesting properties that make them suitable to be used in data mining, data clustering and other application areas.

In the last two decades, many advances in the computer sciences have been based on the observation and emulation of processes of the natural world. The origins of *bioinspired informatics* can be traced to the development of perceptrons and artificial life, which tried to reproduce the mental processes of the brain and biogenesis respectively, in a computer environment [1]. Bioinspired informatics also focuses on observing how the nature solves situations that are similar to engineering problems we face.

With the boom of high-speed networks and increasing storage capacity of database clusters and data warehouses, a huge amount of various data can be stored. *Knowledge discovery* and *Data mining* is not only an important scientific branch, but also an important tool in industry, business and healthcare. These techniques target the problematic of processing huge datasets in reasonable time – a task that is too complex for a human. Therefore computer-aided methods are investigated, optimized and applied, leading to the simplification of the processing of the data. The main goal of computer usage is data reduction preserving the statistical structure (clustering, feature selection), data analysis, classification, data evaluation and transformation.

## 1.2 Ant Algorithms

Ant colonies inspired many researchers to develop a new branch of stochastic algorithms: *ant colony inspired algorithms*. Based on the ant metaphor, algorithms for both static and dynamic combinatorial optimization, continuous optimization and clustering have been proposed. They show many properties similar to the natural ant colonies, however, their advantage lies in incorporating the mechanisms, that allowed the whole colonies to effectively survive during the evolutionary process.

Cemetery formation and brood sorting are two prominent examples of insects' collective behavior. However, other types of ant behavior have been observed, for example predator-prey interaction, prey hunting, etc. The most important are mentioned below.

By replicating the behavior of the insects, the underlying mechanisms may be found and a better understanding of nature may be furthermore achieved. By applying the social insect behavior to computer science, we may achieve more effective techniques. Computer models based on the clustering and sorting of insects can lead to better performance in areas such as search, data mining, and experimental data analysis.

## 1.3 Text Extraction

The accuracy for relation extraction in journal text is typically about 60 % [6]. A perfect accuracy in text mining is nearly impossible due to errors and duplications in the source text. Even when linguists are hired to label text for an automated extractor, the inter-linguist disparity is about 30 %. The best results are obtained via an automated processing supervised by a human [7].

Ontologies have become an important means for structuring knowledge and building knowledge-intensive systems. For this purpose, efforts have been made to facilitate the ontology engineering process, in particular the acquisition of ontologies from texts.

## 1.4 Motivation

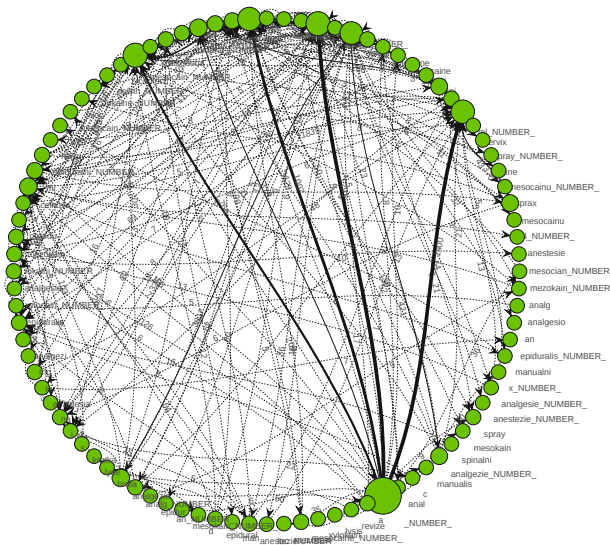
The task of this work is to provide the researchers with a quick automated or semi-automated view on the textual records. Textual data are not easy to

visualize. The word frequency method is simple, but did not provide easily interpretable data. Therefore we decided to extract information in the form of a transition graph.

Such graphs allow us to induce a set of rules for information retrieval. These rules serve for extraction of (boolean/nominal) attributes from the textual rules. These attributes are used in automated rule discovery and can be further used for recommendation. The overall goal of the project is asphyxia prediction during delivery. High asphyxia might lead to several brain damage of the neonate and when predicted, caesarean section might be indicated on time.

## 2 Input Dataset Overview

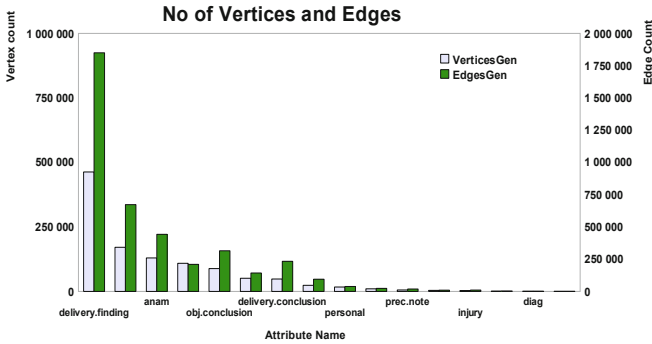
The dataset consists of a set of approx. 50 to 120 thousand records (structured in different relational DB tables; some of them are not input, therefore the range is mentioned)  $\times$  approx. 20 attributes. Each record in an attribute contains about 800 to 1500 characters of text (diagnoses, patient state, anamneses, medications, notes, references to medical stuff, etc.). For textual mining, 16 attributes are suitable (contain sufficiently large corpus).



**Fig. 1.** Figure shows a transitional diagram (directed graph) structure of single attribute literals (a subset). Circular visualization has been used to present the amount of literal transitions (vertices).

The overview of one small (in field length) attribute is visualized in Fig. [1]. Only a subsample (about 5 %) of the dataset could be displayed in this paper, as the whole set would render into an uncomprehensible black stain. The vertices

(literals) are represented as coloured circle, the size reflects the literal (i.e. word) frequency. Edges represent transition states between literals (i.e. the sequence of 2 subsequent words in a sentence/record); edge stroke shows the transition rate (probability) of the edge. The same holds for all figures showing the transition graph, only a different visualization approach has been used.



**Fig. 2.** The DB contains 16 textual attributes that are susceptible for information retrieval via natural language literal extraction. Number of literals (vertices) and transitions (edges) in the probabilistic models are shown for each attribute in a left/right bar respectively. Note the different y-axis scales.

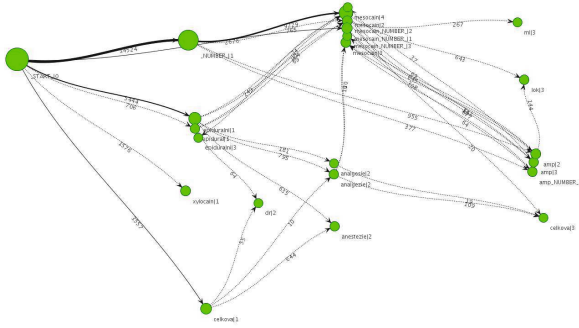
It is clear, that human interpretation and analysis of the textual data is very fatiguing, therefore any computer aid is highly welcome.

## 2.1 Graph Explanation

In this paper we describe *transition graphs*. These are created for each attribute. An attribute consists of many records in form of a sentence. By *sentence* we hereby mean a sequence of literals, not a sentence in a linguistic form. The records are compressed – unnecessary words (such as verbs *is*, *are*) are omitted. In this paper, only the attribute describing the anesthetics during delivery is visualized, as it is the simplest one.

Vertices of the transition graph represent the words (separated by spaces) in the records. For each word (single or multiple occurrence) a vertex is created and its potency (number of occurrences) is noted. For example, the words *mesocaine*, *anesthetics*, *not*, *mL* form a vertex. Note that also words as *mesocain*, *mezokain* and other versions of the word *mesocaine* are present. For a number (i.e. sequence of digits) a special literal *\_NUMBER\_* is used.

Edges are created from single records (sentences entered). For example the sentence *mesocaine 10 mL* would add edges from vertex *mesocaine* to vertex *\_NUMBER\_* and from vertex *\_NUMBER\_* to the vertex *mL* (or the edge count is increased in case it exists). For all records, the count of the edges is also useful.



**Fig. 3.** Very nice graph (sub-graph) providing the basic information about the attribute presented. Note that similar words are clustered (positioned nearby) and the flow of the most common sentences can be easily traced.

It provides an overview on the inherent structure of the data – the most often word transitions.

### 3 Nature Inspired Techniques

Social insects, i. e. ant colonies, show many interesting behavioral aspects, such as self-organization, chain formation, brood sorting, dynamic and combinatorial optimization, etc. The coordination of an ant colony is of local nature, composed mainly of indirect communication through pheromone (also known as *stigmergy*).

The high number of individuals and the decentralized approach to task coordination in the studied species means that ant colonies show a high degree of parallelism, self-organization and fault tolerance. In studying these paradigms, we have high chance to discover inspiration concepts for many successful meta-heuristics.

#### 3.1 Ant Colony Optimization

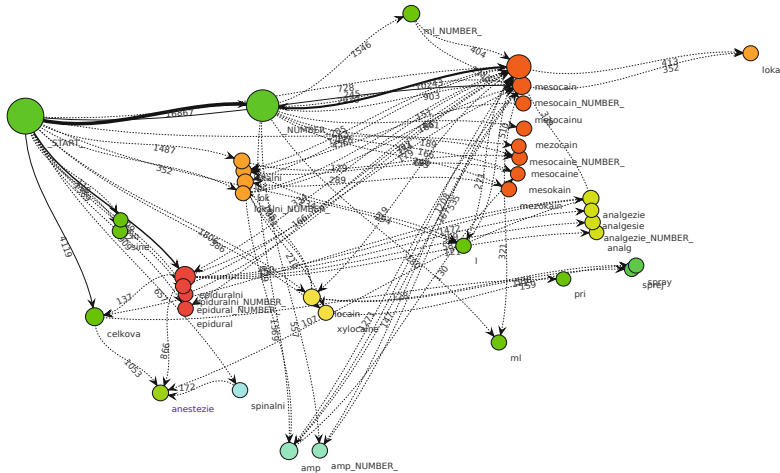
Ant Colony Optimization (ACO) [5] is an optimization technique that is inspired by the foraging behavior of real ant colonies. Originally, the method was introduced for the application to discrete and combinatorial problems.

#### 3.2 Ant Colony Methods for Clustering

Several species of ant workers have been reported to form piles of corpses (ceme-teries) to clean up their nests. This aggregation phenomenon is caused by attraction between dead items mediated by the ant workers.

This approach has been modeled in the work of Deneubourg et al. [4] and in the work of Lumer and Faieta [8] to perform a clustering of data.





**Fig. 5.** An expert (human) organized transition graph (sub-graph) showing the most important relations in one textual attribute. Refer to section [2].

more logical manner. Time needed to organize such graph was about 5–10 minutes. The problem is that the transition graph contains loops, therefore the manual organization is not straightforward.

An aid of a human expert has been used in semi-automated approach (see Fig. [6] where the automated layout has been corrected by the expert. The correction time has been about 20–30 seconds only.

## 6 Parallelization

The ACO\_DTree algorithm has been parallelized in order to take advantage of multicore processors. It contains naturally parallelizable parts, such as population evaluation and population improvement (via the PSO method). Experimental tests have been performed on the 4-core i7-2600 CPU@3.40 GHz (8 cores with Hyperthreading) processor. Performance tests have been run with varying number of cores with and without hyperthreading (HT). The number of execution threads has been increased from 1 to 16.

The *CPU utilization* (load) scaled up to the number of cores (regardless of the HT setting) linearly. There has been a drop-down in CPU load when the no. of threads increased over the number of cores available. The performance for 2 and 4 cores (w/o HT) and 4 and 8 cores (HT) has been similar.





of Biomedical Engineering II” of the CTU in Prague, sponsored by the Ministry of Education, Youth and Sports of the Czech Republic and by the project number NT11124-6/2010 ”Cardiotocography evaluation by means of artificial intelligence” of the Ministry of Health Care. This work has been developed in the BioDat research group `bio.felk.cvut.cz`.

## References

1. Adami, C.: *Introduction to Artificial Life*. Springer (1998)
2. Bursa, M., Huptych, M., Lhotska, L.: Ant colony inspired metaheuristics in biological signal processing: Hybrid ant colony and evolutionary approach. In: *Biosignals 2008-II*, vol. 2, pp. 90–95. INSTICC, Setubal (2008)
3. Bursa, M., Lhotska, L., Macas, M.: Hybridized swarm metaheuristics for evolutionary random forest generation. In: *Proceedings of the 7th International Conference on Hybrid Intelligent Systems 2007 (IEEE CSP)*, pp. 150–155 (2007)
4. Deneubourg, J.L., Goss, S., Franks, N., Sendova-Franks, A., Detrain, C., Chretien, L.: The dynamics of collective sorting robot-like ants and ant-like robots. In: *Proceedings of the first International Conference on Simulation of Adaptive Behavior on From Animals to Animats*, pp. 356–363. MIT Press, Cambridge (1990)
5. Dorigo, M., Stutzle, T.: *Ant Colony Optimization*. MIT Press, Cambridge (2004)
6. Freitag, D., McCallum, A.K.: Information extraction with hmms and shrinkage. In: *Proceedings of the AAAI Workshop on Machine Learning for Information Extraction* (1999)
7. Lafferty, J., McCallum, A., Pereira, F.: Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In: *Proceedings of the ICML*, pp. 282–289 (2001); Text processing: interobserver agreement among linguists at 70
8. Lumer, E.D., Faieta, B.: Diversity and adaptation in populations of clustering ants. In: *From Animals to Animats: Proceedings of the 3th International Conference on the Simulation of Adaptive Behaviour*, vol. 3, pp. 501–508 (1994)