

Advanced Information Technology - Support of Improved Personalized Therapy of Speech Disorders

M. Danubianu, S.G. Pentiuc, I. Tobolcea, O.A. Schipor

Mirela Danubianu, Stefan Gheorghe Pentiuc, Ovidiu Andrei Schipor

“Ștefan cel Mare” University of Suceava
Romania, 720229 Suceava, 13 Universității
E-mail: {mdanub, pentiuc, schipor}@eed.usv.ro

Iolanda Tobolcea

“Alexandru Ioan Cuza” University of Iași
Romania, 700506 Iasi, 11 Bulevardul Carol I
E-mail: itobolcea@yahoo.com

Abstract: One of the key challenges of the Sustainable Development Strategy adopted by the European Council in 2006 is related to public health whose general objective envisages a good level of public health. One of the specific targets includes better treatments of diseases. It is true that there are affections which by their nature do not endanger the life of a person, however they may have a negative impact on her/his life standard. Various language or speech disorders are part of this category, but if they are discovered and treated in due time, they can be often corrected. The difficulty for researchers and therapists is to identify those children who have disorders that show a wide range of issues that cannot be solved spontaneously or which may lead to further significant deficiencies. Information technology in the latest years was used by specialists in order to assist and supervise speech disorder therapy. Consequently they have collected a considerable volume of data about the personal or familial anamnesis, regarding various disorders or regarding the process of personalized therapies. These data can be used in data mining processes that aim to discover interesting patterns which can help the design and adaptation of different therapies in order to obtain the best results in conditions of maximum efficiency. The aim of this paper is to present the Logo-DM system. This is a data mining system that can be associated with TERAPERS system in order to use the data from its database as a source for analysis and to provide new information based on an improved system of therapy. Through the use of appropriate techniques of data mining Logo-DM realizes predictions on the evolution and the final status of patients undergoing therapy and enriches the knowledge data of expert system embedded in TERAPERS.

Keywords: personalized therapy, data mining, classification, clustering, associations rules.

1 Introduction

Various forms of speech disorders affect an important percent of people. There are affections which, by their nature, do not endanger the life of a person, however may have a negative impact on her/his life standard. Discovered and treated in due time, they can be corrected, most often during childhood. The use of information technology in order to assist and supervise speech disorder therapy allows specialists to collect a considerable volume of data about the personal or familial anamnesis, regarding various disorders or regarding the process of personalized therapy.

Even if these data can provide plenty of statistical information little useful knowledge can be obtained from it. In order to get such useful knowledge it is necessary to discover patterns in the data regarding the common characteristics of children with different types of diagnosis, about the connection between antecedents, personal and family behaviour and evolution of the child, or on the connection between the anamnesis and the response to different types of treatments or to different phases of the therapeutic process. These patterns are used to establish such a future strategy so as to maximize the benefits of the therapy and to minimize the costs.

What are the speech disorders? A speech disorder is a problem with fluency, voice, and/or how a person utters speech sounds. Classifying speech into normal and disorder is complex because the statistics points out that only 5% to 10% of the population has a completely normal manner of speaking, all others suffer from one disorder or another. The most common speech disorders are: stuttering, cluttering, voice disorders, dysarthria and speech sound disorders. The speech disorder therapy should begin as soon as possible. Children enrolled in therapy early in their development (younger than 5 years) tend to have better outcomes than those who begin therapy later. During the therapy, speech therapists use a variety of strategies including: oral motor or feeding therapy, articulation therapy and language intervention activities [2]. During the language intervention activities the therapist will interact with a child by playing and talking. He may use pictures, books, objects, or ongoing events to stimulate language development. The therapist may also model correct pronunciation and use repetition exercises to build speech and language skills.

In the area of speech disorders there are some European projects developed as part of the EU Quality of Life and Management of Living Resources program, like: OLP (Ortho-Logo-Paedia) project [8], STAR - Speech Training, Assessment, and Remediation [12] [19], Speechviewer III developed by IBM [11] or ARTUR (Articulation Tutor) [17] [18]. Currently, the priorities at the international level focus on the development of information systems that can provide a personalised therapy. At the national level, little research has been conducted on the therapy of speech impairments [13]. TERAPERS project [1] [2], developed with the financial support granted by the National Agency for Scientific Research, contract ref. no. 56-CEEX-II03/27.07.2006 by the Research Center for Computer Science in the University "Stefan cel Mare" of Suceava, aims to assist and support the speech disorder therapists in their efforts to develop personalized programs for the therapy of dyslalia.

2 Data mining and its application in logopaedic area

Data mining is defined as the process of discovering non-obvious and potentially useful patterns in large data volumes. As exploration and analysis technique of large amounts of data in order to detect patterns or rules with a specific meaning, data mining may facilitate the discovery from apparently unrelated data, relationships that can anticipate future problems or might solve the studied problems.

Data mining represents one phase in the complex process of knowledge discovery in databases (KDD) [5]. According to CRISP-DM [15], the reference model for this process, KDD consists of a sequence of steps. These steps are presented in Figure 1.

Using appropriate methods, data mining can solve two broad categories of problems: prediction and description [10] [14]. The most used methods for prediction are classifications and regressions, and for description, clustering, deviation detection or association rules.

The specific logopaedic tasks performed by data mining fall into the following categories [3]:

- classification which places the people with different speech impairments in predefined classes. Thus it is possible to track the size and structure of various groups. We can use

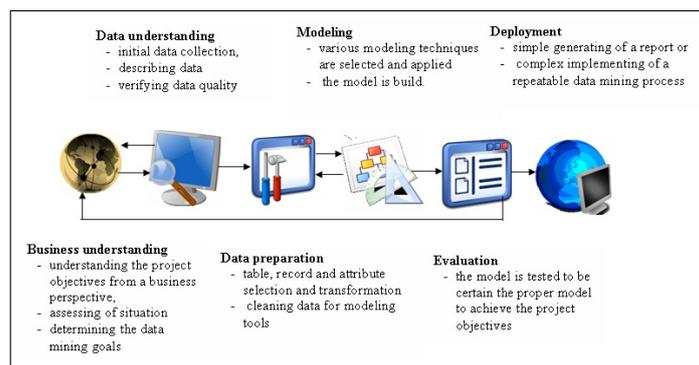


Figure 1: Crisp_DM process of Knowledge Discovery in Databases

classification which is based on the information contained in many predictor variables, such as personal or familial anamnesis data or related to lifestyle, to join the patients with different segments.

- clustering which groups people with speech disorders on the basis of similarity of different features. It is an important task because it helps therapists understand their patients. Clustering aims to finding subsets of a predetermined segment, with homogeneous behavior towards various methods of therapy that can be effectively targeted by a specific therapy but it is not based on the previous definition of groups.
- association rules aim to find out associations between different data which seem to have no semantic dependence. It may be a way to determine why a specific therapy program has been successful on a segment of patients with speech disorders and on the other was ineffective.

To conclude with we state that data mining can be a useful tool. Still, there is a limitation we have to consider. Data mining applications generate information by analyzing patterns of data obtained from the systems which assist and supervise the speech therapy. Such patterns can help predict the evolution of the individuals that are currently in the process of therapy, or design a scheme of an appropriate therapy for them. However data mining technology can not provide information about impairments, people or behaviors that are not found in the databases that provide data for analysis.

3 Logo-DM System

3.1 Objectives

The idea of trying to improve the quality of logopaedic therapy by applying some data mining techniques started from TERAPERS project developed within the Research Center for Computer Science in the University "Stefan cel Mare" of Suceava. This project has proposed to develop a system which is able to assist speech therapists in their speech therapy of dislalya and to assess how the patients respond to various personalized therapy programs. Starting in March 2008 the system is currently used by the therapists from Regional Speech Therapy Center of Suceava.

At present, because of the limited time and the economical aspects involved, information regarding the therapy for each particular case is of interest [4]: what is the predicted final state for a child or what will be his/her state at the end of various stages of therapy, which the best

exercises are for each case and how they can focus their effort to effectively solve these exercises or how the family receptivity - which is an important factor in the success of the therapy - is associated with other aspects of family and personal anamnesis. All this may be the subject of predictions obtained by applying data mining techniques on data collected by using a computer based therapy system. It is also interesting, as part of the knowledge discovered by data mining algorithms, to be used to enrich the knowledge base of expert system embedded. To achieve this goal we propose the development of Logo-DM system.

Consequently its objectives are:

- analysis of data collected and their preprocessing in order to assure a proper quality for data mining algorithms
- feature selection for the elimination of those irrelevant or redundant
- the use of corresponding data mining methods and algorithms that can be applied in order to find models which can answer to problems raised in speech disorders therapy
- models evaluation and their validation on new cases
- to find new rules which can enrich the knowledge base of the expert system embedded in TERAPERS

3.2 System Architecture

Data mining aims at deriving knowledge from data. The architecture of a data mining system plays an important role in the efficiency with which data is mined. Considering the characteristic of the domain we have proposed for the system a two tier client server architecture. This architecture is presented in Figure 2.

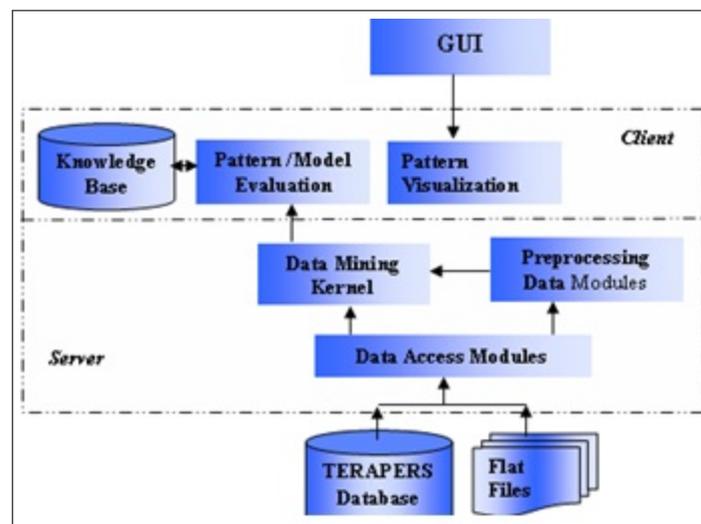


Figure 2: Logo-DM Architecture

On the client side there is the user interface (GUI) which allows the user to communicate with the system in order to select the task to perform, to select and submit the datasets on which data mining needs to be applied. Pattern evaluation and the post-processing step consisting in pattern visualization are performed also on the client. The knowledge base is the module where the background knowledge is stored.

The more difficult computational tasks of data mining operations are carried out on the server. Here, the data mining kernel contains modules able to perform classifications and association rule detection. Supplementary the pre-processing data module allows data to become suitable for applying data mining algorithms.

3.3 Some aspects regarding the system implementation

It is well known that the best results of data mining algorithms are obtained by applying on data in data warehouses. But in this case the development of a data warehouse is not appropriate, so, it is used, as the primary source of data, a database that contains data collected from the different speech therapists' offices. In order to choose the right solution for the implementation of the system we have made an analysis of available data both its structure and content.

We have started from a scheme with over 60 tables and after deleting tables with irrelevant content for the intended purpose we have obtained, as underlying tables for the final data set, 27 tables as presented in Figure 3.

Content analysis can reveal interesting issues related to data quality or the need for transformation. We have made a first assessment of data quality through the following measures: completeness, conformity, accuracy, consistency and redundancy. The mechanisms provided by the used database management system have imposed a minimum, controlled redundancy and have assured data consistency. Values stored in fields correspond to reality, but unfortunately in some records useful data for analysis are missing. Therefore it is necessary to supplement data gaps, and where not possible, the removal of the record for accurate results is suggested.

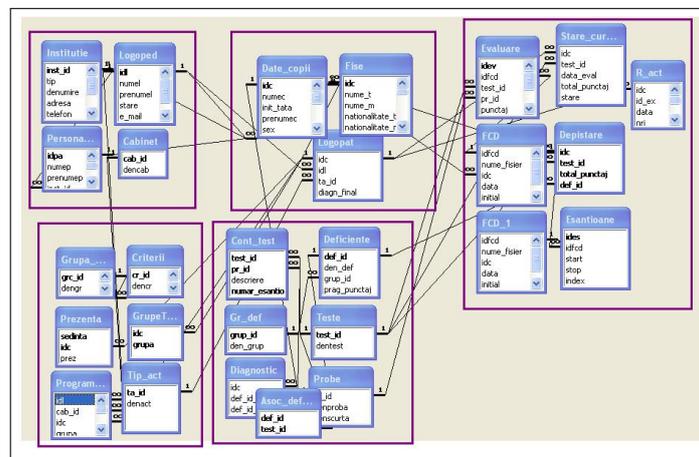


Figure 3: The useful part of database schema

Proper data for the analysis are subjected to the following types of transformation: transformations of the structure, and changes aimed value.

Structural transformations are dictated by the fact that there are fields in the database containing data related to a complex of features to be addressed individually in the analysis. Values of transformations refer to the replacement of coded data by the rules, enabling, for example, the effective storage with descriptive values of characteristics allowing rapid interpretation of results.

An example of these transformations is the following. An issue addressed in the anamnesis form is related to the skills of the child. In Figure 4 we can see that there is a complex of skills of interest (verbal, perceptual, numeric, psycho-motor or special skills).

In the database, all these skills are in two distinct fields: one for general skills, which groups data regarding verbal, perceptual, numeric, psycho-motor and intelligence skills and one for

Figure 4: Sample of anamnesis data

special skills (Figure 5). The field called '*aptitudini*' is numeric and is represented in the table by a string of five bits, as shown in Figure 5. These bits, positioned from left to right, have the following meaning:

- the first bit - verbal skills (1- present, 0- absent)
- the second bit - perceptual skills (1- present, 0- absent)
- the third bit - numeric skills (1- present, 0- absent)
- the fourth bit - psycho-motor skills (1- present, 0- absent)
- the fifth bit - intelligence (1- normal intelligence, 0 - mental deficiency)

emotivitate	disp_afect	aptitudini	apt_spec	atitudini	diagnos
0	<input type="checkbox"/>	0	0	0	
0	<input type="checkbox"/>	0	0	0	
0	<input checked="" type="checkbox"/>	11110	110010	111	
0	<input checked="" type="checkbox"/>	11110	0	111	

Figure 5: Data to be transformed

Since all these attributes may affect the analysis it is desirable that they can be addressed individually and explicitly in the final data set. For this purpose the original table structure is changed and values are converted to descriptive values as in Figure 6.

These changes have conducted to a modified form of the relational database used by Terapers. In the first phase, construction of target data sets for each of the methods to be applied in the system is through the application of relational expressions like those presented in (1).

$$\prod_{I_i} (T_1 \triangleright \triangleleft T_2 \triangleright \triangleleft \dots \triangleright \triangleleft T_k) \quad (1)$$

where:

- I_i is a superset of the attributes regarding the useful characteristics for each method
- $T_1 \dots T_k$ is the set of tables containing the attributes in the list of projection.

apt_verb	apt_erc	apt_num	apt_pm	inteligenta
Prezente	prezente	prezente	prezente	deficenta mintala

aptitudini
0
0
11110
11110

Figure 6: Transformed data

Each of these expressions was implemented in SQL, and has generated intermediate tables. For example, the target data set necessary to establish the profile of children with speech disorders, can be obtained by joining tables which contain: general data about children, family and personal an-amnesis, data on complex evaluation and diagnosis associated. The statement that performs that is presented in (2). The result is a table that contains 129 features.

```

create table caract_copii as
select f.*, l.diagn_final
from fise f, logopat l
where f.idc = l.idc;

```

(2)

Data mining techniques were not designed to process large amounts of irrelevant features. Consequently before their application, a selection of the relevant features is required [6] [7]. The most important objectives of feature selection are: to avoid over fitting and improve model performance. A variant of the mRMR method [9] for categorical values has been used for feature selection. It is based on mutual information criteria, formally defined, for two discrete random variables X and Y , as:

$$I(X;Y) = \sum_{y \in Y} \sum_{x \in X} p(x,y) \log \left(\frac{p(x,y)}{p_1(x)p_2(y)} \right) \quad (3)$$

where $p(x,y)$ is joint probability distribution function of X and Y , and $p_1(x)$ and $p_2(y)$ are the marginal probability distribution functions of X and Y respectively.

For discrete random variable, the joint probability mass function is:

$$p(x,y) = p(X = x, Y = y) = p(Y = y|X = x) * p(X = x) = p(X = x|Y = y) * p(Y = y) \quad (4)$$

Since these are probabilities, we have

$$\sum_x \sum_y p(X = x, Y = y) = 1 \quad (5)$$

The marginal probability function, $p(X = x)$ is:

$$p(X = x) = \sum_y p(X = x, Y = y) = \sum_y p(X = x|Y = y)p(Y = y) \quad (6)$$

The criterion used is related to minimizing redundancy and maximizing relevance to the chosen characteristics. The result of tests performed on data prepared as described in the example mentioned above, revealed that, for classification, the minimum error is obtained if we deal with a number between 20 and 22 features selected. The target data set, obtained after these steps, is

subject to data mining algorithms. For an effective implementation of algorithms we have taken into account, and we tested, two possibilities: to use the Oracle Data Mining kernel (ODM) which offers the possibility to apply algorithms for classification, clustering and association rules and to use some open source implementations of relevant algorithms adapted and integrated into our own system.

We took into account the types of data included in the set and we used implementations in Oracle of Adaptive Bayes Network, Seeker Model and decision trees build with CART [16] and ID3/C4.5 for classification, for clustering the Oracle implementation of A-Clustering algorithm and for association rules Apriori algorithm. It should be noted that for the moment, the volume of data on which work is relatively low, because the system which is the main source of these data is operational for only several months.

4 Conclusions and Future Works

Considering the opportunity of data mining techniques application on data collected in the process of speech therapy, we have concluded that methods such as classification, clustering or as-ociation rules can provide useful information for a more efficient therapy. Consequently, we have designed and we are currently implementing a data mining system that aims to use data provided by TERAPERS system, developed by the Research Center for Computer Science in the University "Stefan cel Mare" of Suceava, in order to achieve an optimized personalized therapy of dyslalia. We have tested the modules for data pre-processing and on target data sets obtained from these modules we have applied more algorithms for detecting the most appropriate solutions for the data mining kernel. At present efforts are directed towards the implementation of evaluation patterns and visualization modules and towards building a user friendly interface.

Bibliography

- [1] M. Danubianu, S.G. Pentiuc, O. Schipor, I. Ungureanu, M. Nestor, Distributed Intelligent System for Personalized Therapy of Speech Disorders, in Proc. of *The Third International Multi-Conference on Computing in the Global Information Technology ICCGI*, July 27- August 01, Athens, Greece, 2008.
- [2] M. Danubianu, S.G. Pentiuc, O. Schipor, M. Nestor, I. Ungurean, D.M. Schipor, TERAPERS - Intelligent Solution for Personalized Therapy of Speech Disorders, *International Journal on Advances in Life Science*, p.26-35, 2009.
- [3] M. Danubianu, T. Socaciu, Does Data Mining Techniques Optimize the Personalized Therapy of Speech Disorders?, *Journal of Applied Computer Science and Mathematics*, p.15-19, 2009
- [4] M. Danubianu, S.G. Pentiuc, T. Socaciu, Towards the Optimized Personalized Therapy of Speech Disorders by Data Mining Techniques, *The Fourth International Multi Conference on Computing in the Global Information Technology ICCGI 2009*, Vol: CD, 23-29 August, Cannes - La Bocca, France, 2009
- [5] F.G. Filip, *Decizii asistate de calculator*, Ed. Tehnica, Bucuresti, 2005
- [6] I. Guyon, A. Elisseeff, An introduction to variable and feature selection. *J. Mach Learn Res.*, 3, p.1157-1182, 2003

- [7] H. Liu, H. Motoda, *Feature Selection for Knowledge Discovery and Data Mining*. Kluwer Academic Publishers, Norwell, MA, 1998
- [8] OLP (Ortho-Logo-Paedia) - Project for Speech Therapy (<http://www.xanthi.ilsp.gr/olp>); W.-K. Chen, *Linear Networks and Systems* (Book style). Belmont, CA: Wadsworth, p. 123-135, 1993
- [9] H. Peng, F. Long, C. Ding, Feature Selection Based on Mutual Information: Criteria of Max-Dependency, Max-Relevance, and Min-Redundancy, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 27, No. 8, p. 1226-1238, 2005
- [10] B. Reiz, L. Csató, Bayesian Network Classifier for Medical Data Analysis. *International Journal of Computers Communications & Control* Vol. 4, p: 65-72, 2009
- [11] Speechviewer III - (<http://www.synapseadaptive.com/edmark/prod/sv3>)
- [12] STAR Speech Training, Assessment, and Remediation (<http://www.asel.udel.edu/speech>)
- [13] Tobolcea, I., *Interventii logo-terapeutice pentru corectarea formelor dislalice la copilul normal*, Editura Spanda, Iasi, 2002.
- [14] P. Wessa, Quality Control of Statistical Learning Environments and Prediction of Learning Outcomes through Reproducible Computing, *International Journal of Computers Communications & Control* Vol. 4, p: 185-197, 2009
- [15] R. Wirth, J. Hipp, CRISP-DM: Towards a standard process model for data mining. *In Proceedings of the 4th International Conference on the Practical Applications of Knowledge Discovery and Data Mining*, pages 29-39, Manchester, UK, 2000
- [16] www.salford-systems.com/
last visited October 2009
- [17] www.speech.kth.se/multimodal/ARTUR/index.html
last visited August 2009
- [18] O. Balter, O. Engwall, A.M. Oster, H. Kjellstrom, Wizard-of-Oz Test of ARTUR - a Computer-Based Speech Training System with Articulation Correction. *Proceedings of the Seventh International ACM SIGACCESS Conference on Computers and Accessibility*, Baltimore, October, 2005, pp.36-43.
- [19] H.T. Bunnell, M.D. Yarrington, B.J. Polikoff, Articulation Training for Young Children, *Proceedings of 6th International Conference on Spoken Language Processing (ICSLP 2000)*, Beijing, China, October 16-20, 2000, vol.4, pp. 85-88.