



**European Journal for Biomedical  
Informatics**

**Volume 2 (2006), Issue 1**

ISSN 1801-5603  
[www.ejbi.org](http://www.ejbi.org)

## Aims and Scope:

The European Journal for Biomedical Informatics reacts on the great European need to share the information in the multilingual and multicultural European area. The journal publishes peer-reviewed papers in English and other languages simultaneously. This opens new possibilities for faster transfer of scientific-research pieces of knowledge of many European countries to a large international community of biomedical researchers, physicians, other health personnel and citizens.

## **Editors and Management:**

**Editor in Chief:** Jana Zvárová, Czech Republic

**Managing Editor** Petra Přečková, Czech Republic

**Sales and Marketing Manager** Libor Seidl, Czech Republic

## **Editorial Board:**

Ammenwerth, Elske	Austria
Blobel, Bernd	Germany
Bobrowski, Leon	Poland
Bureš, Vít	Czech Republic
Degoulet, Patrice	France
Dostállová, Tat'jana	Czech Republic
Eryilmaz, Esat Nadir	Turkey
Hanzlíček, Petr	Czech Republic
Iversen, Irma	Norway
Kern, Josipa	Croatia
Lukosevicius, Arunas	Lithuania
Mansmann, Ulrich	Germany
Martin-Sánchez, Fernando	Spain
Masic, Izet	Bosnia and Herzegovina
Mazura, Ivan	Czech Republic
McCullagh, Paul	United Kingdom
Mihalas, Georgie	Romania
Naszlady, Attila	Hungary
Nykänen, Pirkko	Finland
Paralič, Ján	Slovakia
Pisanelli, Domenico M.	Italy
Sharp, Mary	Ireland
Sousa Pereira, Antonio	Portugal
Valenta, Zdeněk	Czech Republic
Vinarova, Jivka	Bulgaria
de Lusignan, Simon	United Kingdom

## **Publisher:**

EuroMISE s.r.o., Paprsková 330/15, CZ-14000 Praha 4, Czech Republic  
EU VAT ID: CZ25666011

Office: EuroMISE s.r.o., Pod Višňovkou 23, CZ-14000 Praha 4, Czech Republic

Contact: Karel Zvára, [zvara@euromise.com](mailto:zvara@euromise.com),

Tel: +420 226 228 904, Fax: +420 241 712 990

## ORGANIZATION OF THE MANUSCRIPT

**Title page.** The first (title) page should contain the title of the paper, names and workplaces of all authors. Individual workplaces are necessary to be graphically differentiated (preferably by numeral as the upper index).

**Abstracts and keywords.** At the beginning the author puts an abstract and keywords. The abstract should be in the extent of 250-300 words. There should be 4 to 7 keywords, according to author's consideration, preferably from MeSH index.

**Main text of the paper.** General rules for writing manuscripts recommend use of simple and declarative sentences; avoid long sentences, in which meaning may be lost by complicated construction. All acronyms and abbreviations should be explained when they first appear in the text. The main text of the paper should follow the style of selected type of paper.

**Acknowledgement.** Acknowledgements, if any, should be given at the end of the paper, before bibliographic references.

**References.** References should be cited in the text by their index number according to the order of appearance in the manuscript. Each reference should be marked by its index number in square bracket corresponding to bibliography section. It is possible to include references to dissertation works and technical reports. It is obligatory to include information sufficient to look up referenced text. Examples of references in bibliography section:

- [1] Knaup P., Ammenwerth E., Brandner R., Brigl B., Fischer G., Garde S., Lang E., Pilgram R., Ruderich F., Singer R., Wolff A. C., Haux R., Kulikowski C.: Towards Clinical Bioinformatics: Advancing Genomic Medicine with Informatics Methods and Tools. *Methods Inf Med* 2004; 43, pp. 302-307
- [2] Blobel B., Pharow P.: A Model-Driven Approach for the German Health Telematics Architectural Framework and the Related Security Infrastructure. In: Connecting Medical Informatics and Bio-Informatics. Proceedings of MIE2005 (Eds. R. Engelbrecht, A. Geissbuhler, C. Lovis, G. Mihalas), Vol. 116, Amsterdam, IOS Press, 2005, pp. 391-396
- [3] <http://www.infobiomed.org/>

**Tables and Figures.** Authors should use tables only to achieve concise presentation, or where the information cannot be given satisfactory in another way. Tables should be numbered consecutively using Arabic numerals and should be referred to in the text by numbers. Each table should have an explanatory caption that should be as concise as possible. Figures should be clear, easy to read and of a good quality. Styles and fonts should match those in the main body of the paper. All figures must be mentioned in the text in consecutive order and should be numbered with Arabic numerals. Authors should indicate precisely in the main text where tables and figures should be inserted, if these elements are given only separately or at the end in the original version of the manuscript.

# Content

<b>José Ignacio Serrano, Marie Tomečková, Jana Zvárová</b>	
<i>Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis</i>	
.....	6
<i>Métodos de Aprendizaje Automático para el Descubrimiento de Conocimiento en Datos Médicos sobre Arterosclerosis (Español).....</i>	34
<i>Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze (Česky).....</i>	62
<b>Daniel Schwarz, Ivo Provazník</b>	
<i>Low-dimensional Multimodal Deformable Registration of MRI Brain Images in Stereotaxic Space.....</i>	90
<i>Málorozměrná multimodální pružná registrace obrazů mozku z MRI ve stereotaktickém prostoru (Česky).....</i>	98
<b>Jana Vrbková, Vilém Bruk</b>	
<i>A Prediction of Blood Flow through a Bypass Graft Using Statistical Methods.....</i>	106
<i>Předpověď průtoku krve bypassem pomocí statistických metod (Česky).....</i>	127

## ORGANIZATION OF THE MANUSCRIPT

**Title page.** The first (title) page should contain the title of the paper, names and workplaces of all authors. Individual workplaces are necessary to be graphically differentiated (preferably by numeral as the upper index).

**Abstracts and keywords.** At the beginning the author puts an abstract and keywords. The abstract should be in the extent of 250-300 words. There should be 4 to 7 keywords, according to author's consideration, preferably from MeSH index.

**Main text of the paper.** General rules for writing manuscripts recommend use of simple and declarative sentences; avoid long sentences, in which meaning may be lost by complicated construction. All acronyms and abbreviations should be explained when they first appear in the text. The main text of the paper should follow the style of selected type of paper.

**Acknowledgement.** Acknowledgements, if any, should be given at the end of the paper, before bibliographic references.

**References.** References should be cited in the text by their index number according to the order of appearance in the manuscript. Each reference should be marked by its index number in square bracket corresponding to bibliography section. It is possible to include references to dissertation works and technical reports. It is obligatory to include information sufficient to look up referenced text. Examples of references in bibliography section:

- [1] Knaup P., Ammenwerth E., Brandner R., Brigl B., Fischer G., Garde S., Lang E., Pilgram R., Ruderich F., Singer R., Wolff A. C., Haux R., Kulikowski C.: Towards Clinical Bioinformatics: Advancing Genomic Medicine with Informatics Methods and Tools. *Methods Inf Med* 2004; 43, pp. 302-307
- [2] Blobel B., Pharow P.: A Model-Driven Approach for the German Health Telematics Architectural Framework and the Related Security Infrastructure. In: Connecting Medical Informatics and Bio-Informatics. Proceedings of MIE2005 (Eds. R. Engelbrecht, A. Geissbuhler, C. Lovis, G. Mihalas), Vol. 116, Amsterdam, IOS Press, 2005, pp. 391-396
- [3] <http://www.infobiomed.org/>

**Tables and Figures.** Authors should use tables only to achieve concise presentation, or where the information cannot be given satisfactory in another way. Tables should be numbered consecutively using Arabic numerals and should be referred to in the text by numbers. Each table should have an explanatory caption that should be as concise as possible. Figures should be clear, easy to read and of a good quality. Styles and fonts should match those in the main body of the paper. All figures must be mentioned in the text in consecutive order and should be numbered with Arabic numerals. Authors should indicate precisely in the main text where tables and figures should be inserted, if these elements are given only separately or at the end in the original version of the manuscript.

# Content

<b>José Ignacio Serrano, Marie Tomečková, Jana Zvárová</b>	
<i>Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis</i>	
.....	6
<i>Métodos de Aprendizaje Automático para el Descubrimiento de Conocimiento en Datos Médicos sobre Arterosclerosis (Español).....</i>	34
<i>Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze (Česky).....</i>	62
<b>Daniel Schwarz, Ivo Provazník</b>	
<i>Low-dimensional Multimodal Deformable Registration of MRI Brain Images in Stereotaxic Space.....</i>	90
<i>Málorozměrná multimodální pružná registrace obrazů mozku z MRI ve stereotaktickém prostoru (Česky).....</i>	98
<b>Jana Vrbková, Vilém Bruk</b>	
<i>A Prediction of Blood Flow through a Bypass Graft Using Statistical Methods.....</i>	106
<i>Předpověď průtoku krve bypassem pomocí statistických metod (Česky).....</i>	127

# Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis

José Ignacio Serrano<sup>1</sup>, Marie Tomečková<sup>2</sup>, Jana Zvárová<sup>2</sup>

*1. Instituto de Automática Industrial, CSIC, Madrid, Spain,*

*2. Department of Medical Informatics, Institute of Computer Science AS CR, Prague, Czech Republic*

Machine learning techniques are methods that given a training set of examples infer a model for the categories of the data, so that new (unknown) examples could be assigned to one or more categories by pattern matching within the model. The data from follow-up studies with repeated collection of the same type of data are very suitable for this analysis. Machine learning algorithms belonging to a variety of paradigms have been applied to knowledge discovery on medical data. All the used algorithms belong to the supervised learning paradigm. Several algorithms have been tested, trying to cover most of the kinds of supervised learning. Two kinds of experiments have been carried out. The first is intended to discover associations between attributes. The second kind is intended to test prediction of future disorders. For the experiments in this paper the data used was from the twenty years lasting primary preventive longitudinal study of the risk factors (RF) of atherosclerosis in middle aged men. Study is named STULONG (LONGitudinal STudy). The results show that some methods predict some disorders better than others, so it is interesting to use all the algorithms at a time and consider the result confidence based upon the known tendency of each method. The machine learning algorithms have been also used in the prediction of death cause, obtaining poor results in this case, maybe due to the small amount of information (entries) of this type in the dataset.

**Keywords:** knowledge discovery, supervised machine learning, biomedical data mining, risk factors of atherosclerosis

## 1. Introduction

Machine learning techniques [1] are methods that given a training set of examples infer a model for the categories of the data, so that new (unknown) examples could be assigned to one or more categories by pattern matching within the model.

Machine learning techniques have been applied successfully to a high variety of problems and data for prediction tasks. The main objective is to research how to apply machine learning algorithms to this data in order to discover relationships between attributes and to make predictions that could be useful for decision support. Medical data is a special kind of data, because many different kinds of features are involved in the collections. Moreover, the medical data have several known problems: missing, incorrect and sparse information and temporal data. Machine learning methods are very suitable for this kind of data [2]. There

exists several KDD works attempting to deal with large-scale medical information. In [3], authors try to detect type of hepatitis by extracting short sequential patterns from the temporal features. In [4], simple rules are discovered using 4ft-miner (i.e. statistical tables of two rows and two columns), in order to temporally characterize, by differences, the hepatitis types B and C. Authors in [5] attempt to discover rules of singles boolean features that can be able to predict the liver fibrosis stage. The same application appears in [6] but, in this case, extracted patterns are clustered and then these clusters are assigned to fibrosis stages depending on the covered examples. They also applied this technique to atherosclerosis risk [7]. The data from follow-up studies with the repeated collection of the same type of data are very suitable for this analysis. Additional examples of data mining on biomedical data are presented in [8] and [9]. For the experiments in this paper the data used was from the twenty years lasting primary preventive longitudinal study of the risk factors (RF) of atherosclerosis in middle aged men. The study is named STULONG (LONGitudinal STUdy) [10]. The main target of this study is to validate machine learning as a way of association mining and to validate classification performance as a measurement of the salience for the discovered association. It is also intended to test machine learning algorithms in the prediction of far future disorders.

In the next section, the details of the STULONG dataset are presented. In section 3, the machine learning algorithms tested are described. Section 4 presents the measures for evaluation and Section 5 describes validation experiments. Finally, concluding remarks and future work is presented in Section 6.

## 2. Description of the Study and Data Set

The STULONG (<http://euromise.vse.cz/challenge2004/index.html>) [10], [11] data were collected by the 2<sup>nd</sup> Department of Internal Medicine, 1<sup>st</sup> Faculty of Medicine of Charles University in Prague and General University Hospital, Prague and transferred to the electronic form and analysed by statistical methods by the European Centre of Medical Informatics, Statistics and Epidemiology of Charles University in Prague and the Academy of Science of the Czech Republic, Prague.

The main aims of the study were:

1. To identify the prevalence of risk factors (RF) of atherosclerosis in a population generally considered as the most endangered by possible complications of atherosclerosis, i.e. middle aged men.
2. To follow the development of these RF and their impact on the examined men health, especially with respect to atherosclerotic cardiovascular diseases.
3. To study the impact of complex intervention of RF on their development and cardiovascular morbidity and mortality in men.

Men born in 1926–1937 and living in Prague 2 were selected from the Prague 2 election lists in year 1975. For the first examination, 1419 of 2370 invited men came. Entry examinations were performed in the years 1976–1979. The invitation for examination included a short explanation of the aims of the study, of the first examination purpose, procedure and later observations and asked for co-operation. At that time, no informed signature of the respondent was required. Should the man react to the first invitation for the examination, we considered that a sufficient agreement with the examination itself, observation and results

Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis

processing. Should man fail to react to the first invitation, we would send two more invitations, minimally.

The risk factors were defined according to the definitions at that time as follows:

- hypertension – blood pressure BP  $\geq 160/95$  mm Hg or men under the hypotensive medication,
- hypercholesterolemia – cholesterol  $\geq 260\text{mg\%}$  (6,7 mmol/l),
- hypertriglyceridemia – triglycerides  $\geq 200\text{mg\%}$  (2,2 mmol/l),
- smoking:  $\geq 15$  cig./day currently or smoking of the same number of cigarettes within 1 year prior to the beginning of the study (pipe or cigar smokers were considered non-smokers),
- overweight: Brocka index  $> 115\%$  (Brocka index: height in cm minus 100 = 100 %),
- positive family case history: death of father or mother from coronary artery disease, or vascular stroke before reaching 65 years of age.

Men were divided according to presence of risk factors (RF), overall health conditions and ECG result into following groups:

- NG** = group of men without RF defined above, without manifestation of the atherosclerotic diseases or other serious illnesses making their ten-year-long observation impossible, and without ECG changes,
- RG** = group of men with at least one RF defined above, without manifestation of the atherosclerotic diseases or other serious illnesses making their ten-year-long observation impossible, and without ECG changes,
- PG** = group of men with a manifested cardio-vascular atherosclerotic diseases or other serious diseases making their ten-year-long observation impossible (e.g. malignant illness, advanced failure of liver or kidneys, serious neurological or psychological problem). The pathologic group included also men with diabetes treated with orally administered anti-diabetics or insulin, and men with pathologic ECG, according to the Minnesota ECG code.

Long-term observation of patients was based on their division into the groups stated above:

- The risk group **RG** was randomly divided into two sub-groups designated as **RGI** (intervened risk group) and **RGC** (control risk group). The patients in the **RGI** group were invited for check up minimally twice a year. Following pharmacological intervention, they were invited as necessary. The patients in the **RGC** group received a short written notice including their laboratory results and ECG description and a recommendation to take these results to their physician; possible intervention of RF was left to the decision of these physicians. At the first examination, no significant difference in age, socio-economic factors or RF occurrence was demonstrated between the RGI and RGC groups.
- 10 % of men in the **NG** group was examined minimally once a year just as the risk group – (they are denoted **NGS**); In this group of men, similarly to the risk group, intervention was initiated as soon as a RF was identified and confirmed (hyperlipidemia, arterial hypertension). The remaining men of the **NG** were invited for a control check up 10–12 years later.
- The men from the **PG** group were excluded from further observation.

Intervention was the key problem of the study and was based on non-pharmacological influence. We tried to modify and to optimize RF.

- *Non-pharmacological intervention:* interviews on lifestyle, i.e. diet, physical activity, suitability or necessity to stop smoking and reduce weight. The interviews were repeated during each control and except for general instructions, they focused also on specific RF of a given man.
- *Pharmacological intervention:* treatment of arterial hypertension and hyperlipoproteinemia – was very limited in the initial stages of the study and may be mostly used only in the last years of the study. The pharmacological therapy was recommended with respect to the overall risk of a given man and his possible other diseases.

Four data files have been used for the analysis:

1. The file *ENTRY* contains values of 244 attributes obtained from entry examinations for each man; these attributes are either codes or results of size measurements of different variables or results of their transformations (identification of man, family and personal history, social factors – education, physical activities, smoking, eating habits, alcohol, after them anthropometric measurements – height, weight, skin folds, physical examination with measurement of blood pressure, pulse, laboratory values and coding of ECG).
2. The file *CONTROL* contains results of observation of 66 attributes recorded during control examinations. There are attributes corresponding to identification, to habit changes, to personal history, physical examination and biochemical values, and data about hypertension, hypercholesterolemia, hypertriglyceridemia and other coronary and oncological diseases. This file consists of 10,572 records of long observation.
3. Additional information about health status of 403 men dropped out in the time of the study was collected by the postal questionnaire. Resulting values of 62 attributes are stored in the *LETTER* file.
4. There are 5 attributes concerning death of 389 patients. Values of these attributes are stored in the *DEATH* file. It contains attributes for the identification of the patients and the date and cause of death.

### 3. Description of the Used Methods

All the used algorithms belong to the supervised learning paradigm. That is, a learning stage is needed in order to build a model over the training examples and then use this model to predict the category of unknown examples. Several algorithms have been tested, trying to cover most of the kinds of supervised learning. Each of the used methods is very briefly explained next:

#### 3.1 Naive Bayes

Naïve Bayes [12] calculates, for each pair attribute-value, for example (*education, university*), the probability of belonging to each category, by dividing the number of examples of the target category where the pair appears by the total number of examples where the pair appears. Thus, each pair will have a probability for each tentative category. Naive Bayes is

based on the assumption that every pair attribute-value within an example is independent on each other. Thus, when an unlabeled example is classified, the probability for each category of the example is the multiplication of the probability for the corresponding category of each of the pairs that form the example. The predicted category is the one with the highest probability.

### 3.2 Multilayer Perceptron

The classification model of the Multilayer Perceptron Neural Network [13] is composed of a certain number of layers of neurons interconnected between them. The architecture used for this dataset is presented in Figure 1.

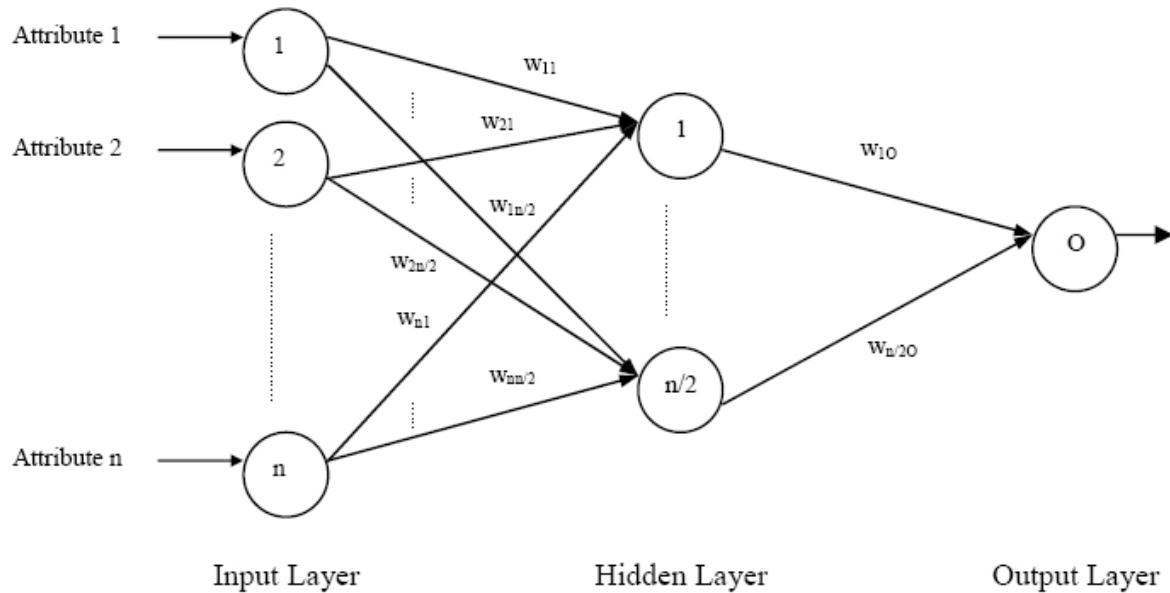


Fig. 1. Architecture of the Multilayer Perceptron Neural Network used.

Each connection has an associated weight. The input to each neuron is the weighted sum, using the association weights, of all the incoming values. The output of each neuron is the result of applying a function. In this case, a typical sigmoid function is implemented in all the neurons. Figure 2 shows the function expression and representation.

$$f(x) = \frac{1}{1 + e^{-x}}$$

Fig. 2. Expression and representation of sigmoid function.

Thus, each of the attribute values from a sample of the dataset is entered in the corresponding neuron of the input layer, and the values spread through the network to the output layer, where the output value of the neuron is the predicted class.

The training phase consists of, given a set of initial weights values, entering each of the labelled examples of the training dataset into the model and comparing the output value with the expected class. Depending on the error of the predicted class, the back propagation algorithm changes weights from the output layer to the input layer, in order to make the predicted value to be more similar to the expected one. This process is carried out a certain number of epochs or iterations. In this case, this number is equal to 500. The amount the weights are changed in back propagation, so called learning rate, is 0.3, and the momentum applied to the weights during updating is 0.2. If the back propagation algorithm does not reach a good approximation to the expected output after one iteration, then it resets the model and causes the learning rate to decrease.

### 3.3 Support Vector Machines (SVM)

Support Vector Machines [14] try to separate examples, based on their category, in the  $n$ -dimensional space, being  $n$  the number of attributes or features, by hyper planes of the form  $\mathbf{w} + \mathbf{b}$ , so that

$$\mathbf{x} \cdot \mathbf{w} + b \geq +1 \rightarrow \text{category} = \text{true}$$

$$\mathbf{x} \cdot \mathbf{w} + b \leq -1 \rightarrow \text{category} = \text{false}$$

$\mathbf{x}$  being the example represented as a vector of  $n$  components. Here,  $\mathbf{w}$  is the support vector, perpendicular to the hyper plane, and correspond to examples that are beyond or over the limits of their category (see Figure 3).

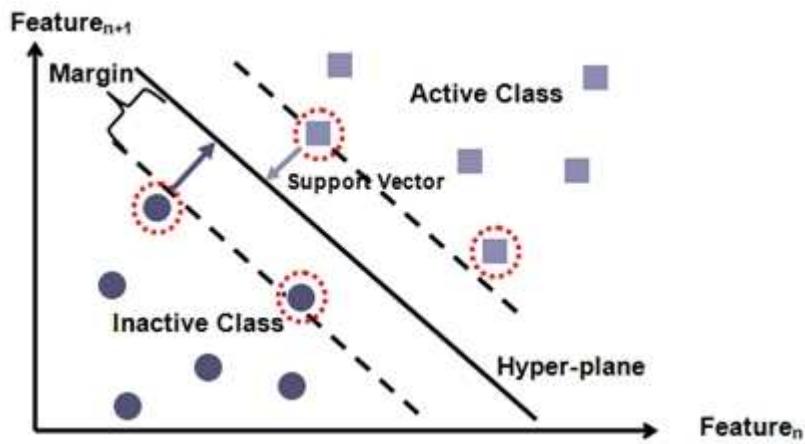


Fig. 3. Support vectors scheme.

The support vector also defines, by its module, a margin of one between the hyper plane and the first positive and negative examples (that is the reason for +1, -1 thresholds). For each category, the algorithm tries to find  $\mathbf{w}$  maximizing the margin. To classify an unlabeled example the algorithm simply applies the expression above. This is a simple implementation

of the method and the one used in the experiments, but there are other more sophisticated implementations and techniques.

### 3.4 K-Nearest Neighbour

KNN is a memory-based algorithm [15], with the background idea that past experiences can help us to solve present ones by analogy. It considers each example as a vector of  $n$  components, being  $n$  the number of attributes or features. It does not need a learning stage. To predict the class of an unlabeled example, the algorithm compares the input example with each of the examples in the training data or memory, by calculating the distance between them. Then, the majority class of the  $K$  most similar training examples is the one predicted for the input example. The distance used in the experiments is the Euclidean distance between vectors. However, there are more possibilities in the literature.

### 3.5 ID3 and C4.5 Decision Trees

The model produced by this algorithm is a tree [16], where each node corresponds to an attribute and each arc of the node corresponds to a possible value of the node attribute.

The learning algorithm constructs the tree from the training data. The selection of the attribute that will form a node, at each moment, is carried out by calculating the entropy of the data after the selection of the node. That is, for each attribute, the entropy of the remaining data without the attribute, separated by the different values of the node attribute, is calculated. Thus, the attribute that produces the minimum entropy is the selected for the node. The process goes on until there is no more attributes or the number of remaining examples under a node is lower than a certain threshold. In the former case, the majority class of these remaining examples is the one settled under the node. In Figure 4, we can see an example:

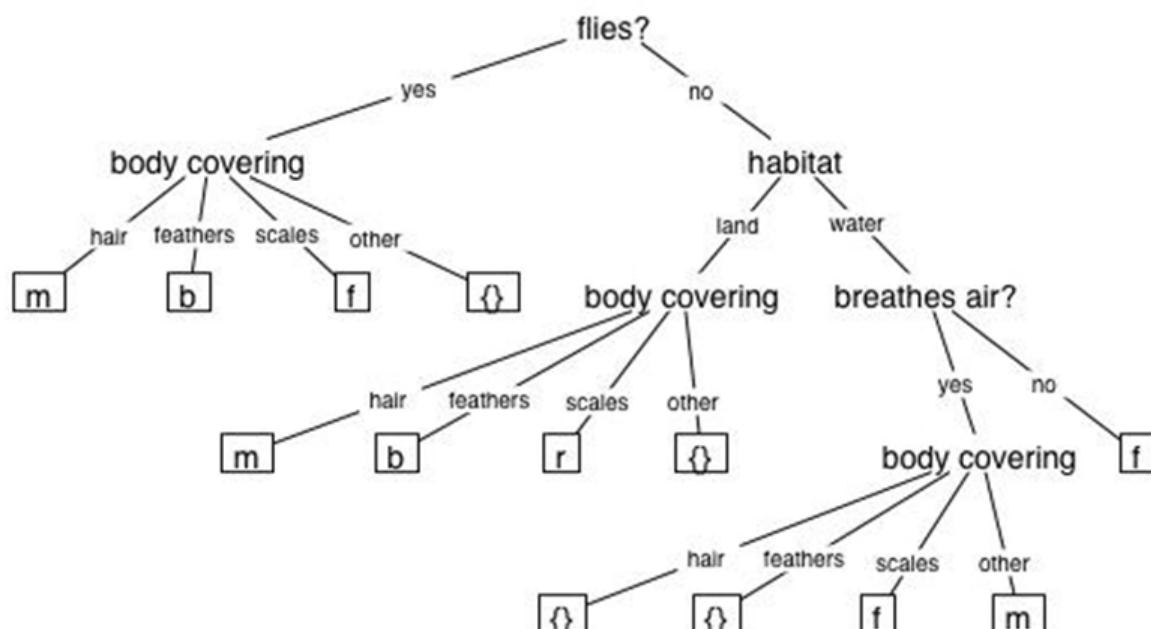


Fig. 4. Example of a decision tree.

In the example there are 4 attributes: *flies*, *body covering*, *habitat* and *breathes air*, and four possible categories, *m*, *b*, *f* and *r*. Here, the first attribute is *flies* because it is the one that produces the division on the data with minimum entropy at that level, and so on. To classify an unlabeled example you only have to follow the tree top-down, and the final leaf is the predicted category. The pathways from the root node to the leaf node can be viewed as rules, where the condition is formed by AND operation of the terms (*node=arc*).

C4.5 is an extension of ID3 that allows continuous numerical attributes, accounts for missing values and carries out a pruning process in order to reduce the tree size for dealing with larger amount of data. The J48 tree used in the experiments is an implementation of C4.5.

### 3.6 Ridor Rules Learner

Ridor stands for the RIpple-DOWN Rule learner [17]. It generates the default rule first and then the exceptions for the default rule with the least (weighted) error rate when it is used to classify the training data. Then it generates the "best" exceptions for each exception and iterates until pure. Thus it performs a tree-like expansion of exceptions and the leaf has only the default rule but no exceptions. The exceptions are a set of rules that predict the classes other than class in the default rule. IREP is used to find out the exceptions. IREP algorithm constructs rules by gradually adding one term in the condition at a time so that the error rate is minimized. The rule condition terms are like (*attribute {=, ≠, ≤, ≥} value*).

## 4. Evaluation

The evaluation processes and measures are the same for all the experiments. Given the data, a part of the collection is considered as a training set and the remaining as a test set. So the models learn from the training set and try to predict the values of examples in the test set. Since the category of test set examples is known, we can check the predictions. Three different typical measures are calculated for each category: precision, recall and F-measure [18]. Precision is the percentage of predictions of one category that were correct. Equation 1 presents the precision expression.

$$\text{Precision}(\text{category}_i) = \frac{\text{number of correct predictions as category}_i}{\text{total number of predictions as category}_i} \quad (1)$$

Recall is the percentage of all the examples of the test set belonging to a category that were correctly predicted. The expression is presented in Equation 2.

$$\text{Recall}(\text{category}_i) = \frac{\text{number of predictions as category}_i}{\text{total number of examples of category}_i} \quad (2)$$

F-measure is a combination of the former measures. It accounts, someway, the intersection between the examples involved in precision and recall, normalized by the sum of both. Equation 3 shows the F-measure expression.

$$F\text{measure} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3)$$

Thus, these three measures are calculated for each category of the test set. As said before, given the collection it is needed to divide the data in training and test sets. A common way of evaluation is cross-validation. The collection is divided into  $n$  equally sized parts. Then, each  $n-1$  part combination is considered as training and the remaining part as test, so the algorithm is run  $n-1$  times, and the final results are the average of this  $n-1$  executions. For all the experiments described below,  $n$  has a value of 3, so training is always 66 % of the data and test stands for 33 %, running each algorithm three times. Usually, the value of  $n$  is greater than 3, typically of 10, but in this case we have very few examples of some categories, and a greater value of  $n$  could produce test sets with no representation of the mentioned categories, what is not desirable.

## 5. Experiments

Two kinds of experiments have been carried out. The first is intended to discover associations between attributes by considering the classification performance as an indicator of the association strength. The second kind is intended to test the prediction of future disorders.

It is needed to remark that the observations in the data set with missing values were not removed nor imputed, because the implementations of the learning algorithms are able to deal with missing data. These implementations are the ones included in the WEKA environment [19], used with default parameters to perform the experiments above.

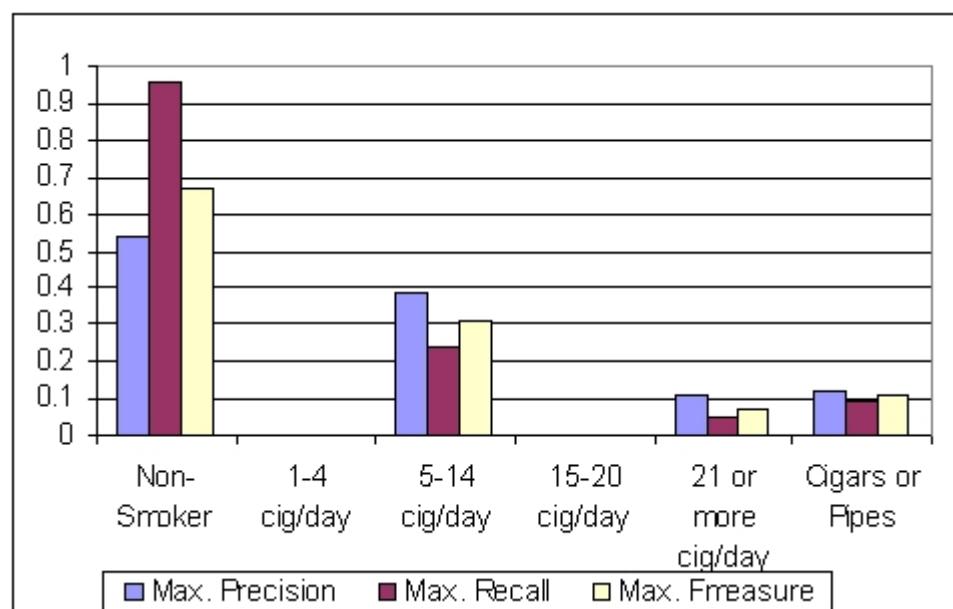
### 5.1 Finding Answers

The first experiments are related to the analytical questions proposed for the Discovery Challenge of ECML/PKDD 2004 conference, specifically the ones related to the Entry collection. These tasks consist of finding relations in three different groups of patients: normal group, risk group and pathologic group. These groups correspond to the risk level of atherosclerosis – see above, and will be referenced as level groups. Specifically, the target relationships are between social factor features and physical activities features, alcohol features, smoking features, body mass index, blood pressure and HDL cholesterol, and then between physical activities and the remaining and between alcohol and the remaining. So, machine learning algorithms are applied to the data of each different group, trying to predict the value of each of the features of one group given the features of the other, viewing the possible values as the considered categories. Thus, for example, given the four social factor attributes as training factors the algorithms are run in order to predict the value of each of the four physical activities attributes and so on with the other features groups. For each relationship, the maximum values over all the different algorithms results are calculated in order to compare between level groups. So, if the prediction accuracy is good, we could say

Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis

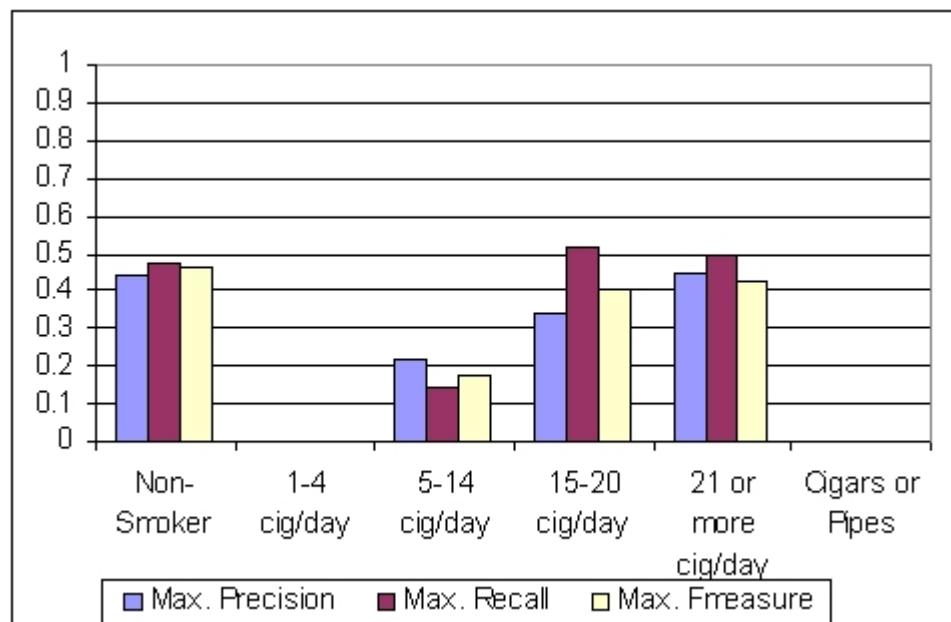
that there is a strong relationship, in a degree equal to the accuracy, between the features used for training and the feature whose values are predicted, and also we can compare prediction measures between features and level groups to state which relations are stronger than others.

Due to paper length limitations, only some of the most representative results are presented. In Figure 5, the maximum precision, recall and F-measure predictions results for "Smoking", given social attributes, are showed for each of the level groups, a) Normal, b) Pathologic and c) Risk, respectively, and given physical activity attributes for each of the level groups, d) Normal, e) Pathologic and f) Risk, respectively. As can be seen, in the Normal group, either from social factors or physical activity, the best prediction is reached for non-smoking people, being not significant for the remaining values of the "Smoking" attribute. It seems that the relationship between social factors and smoking is slightly stronger than physical activity and smoking, because it produces better results for all the values of the "Smoking" attribute. In both Pathologic and Risk groups, the relationship between the training factors and the non-smoking value is stronger for physical activity factors, being in particular high in the Pathologic group. In the latter groups, people who smoke 15 or more cigarettes a day are better predicted than in the Normal group but non-smokers are much worse detected than in the Normal group.

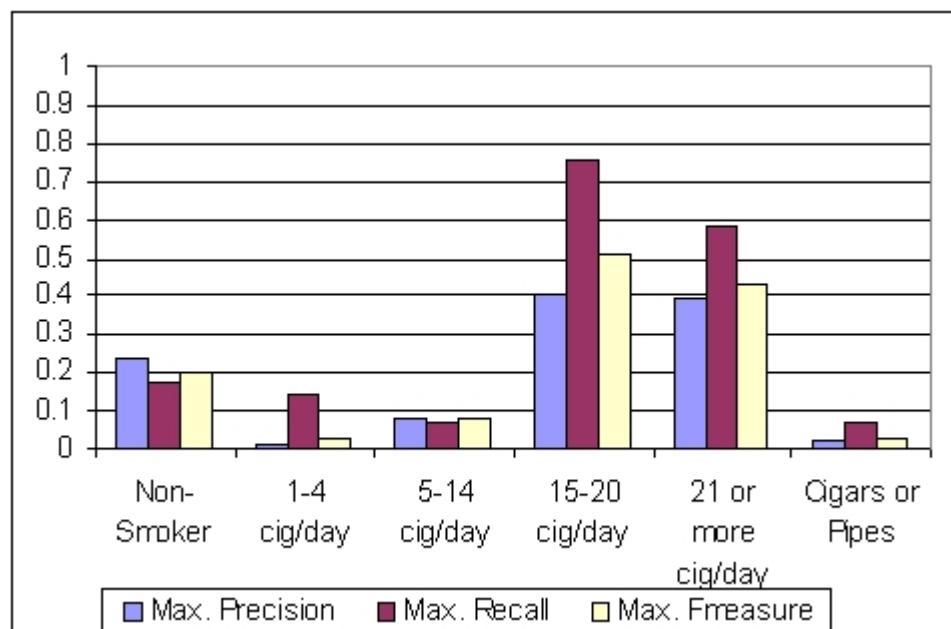


a)

Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis

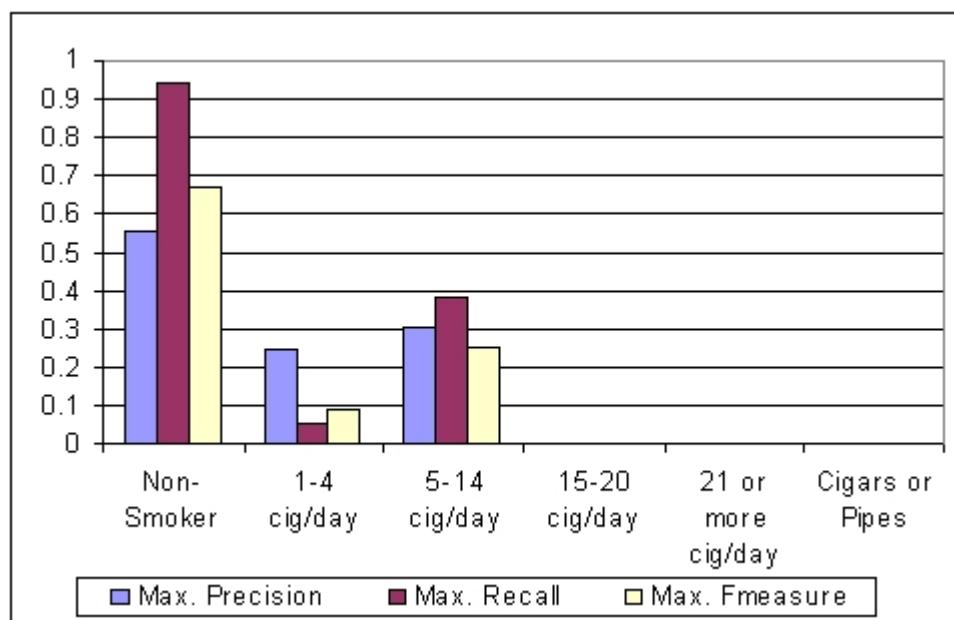


b)

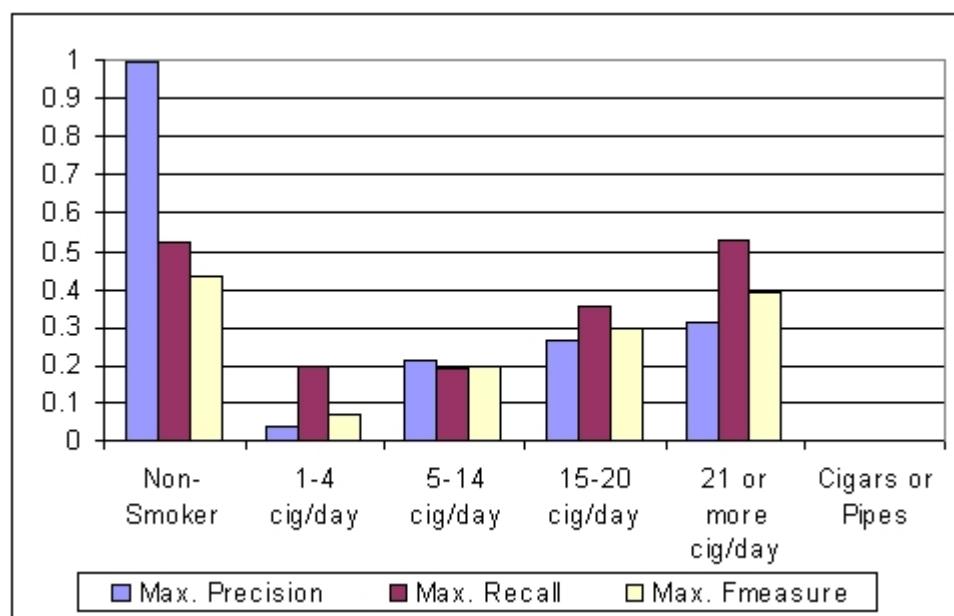


c)

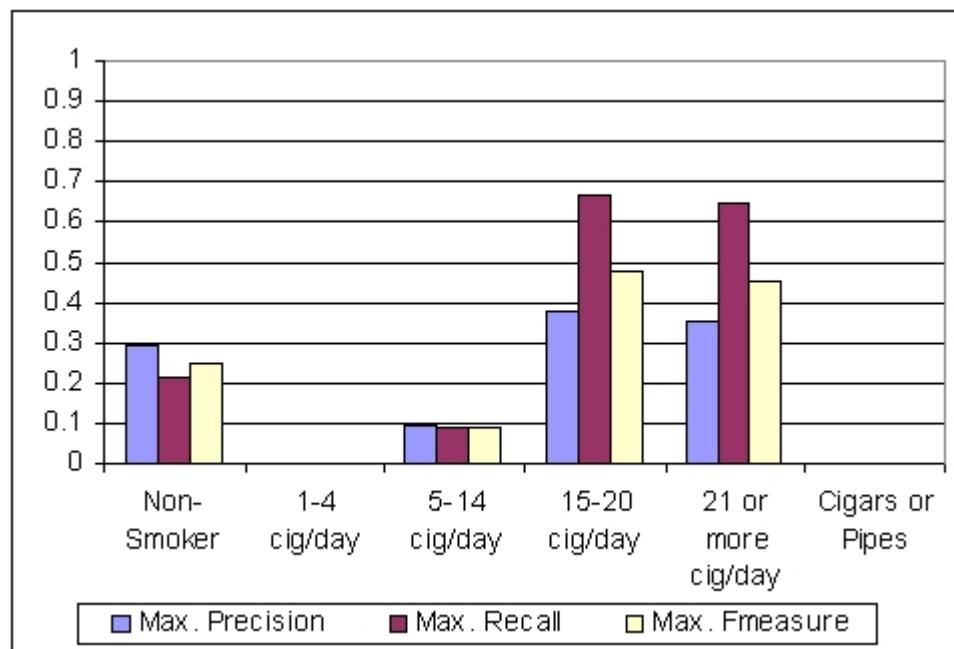
Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis



d)



e)

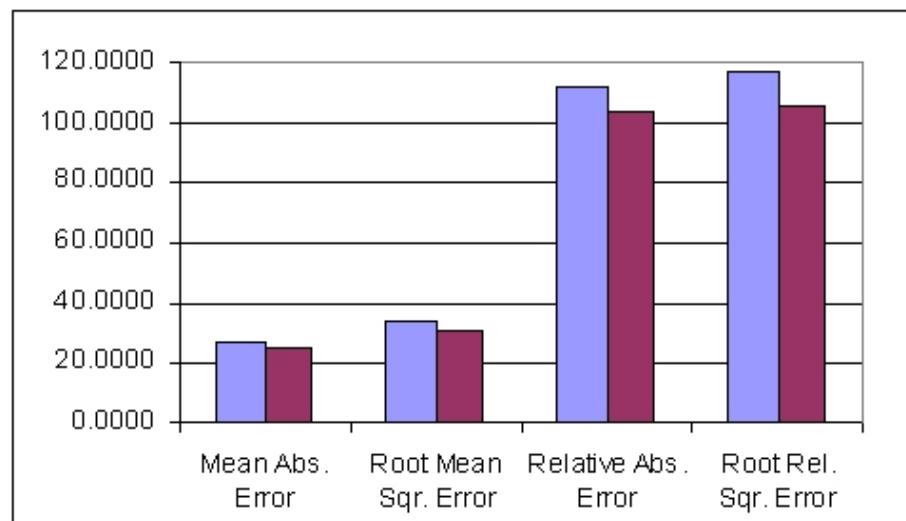


f)

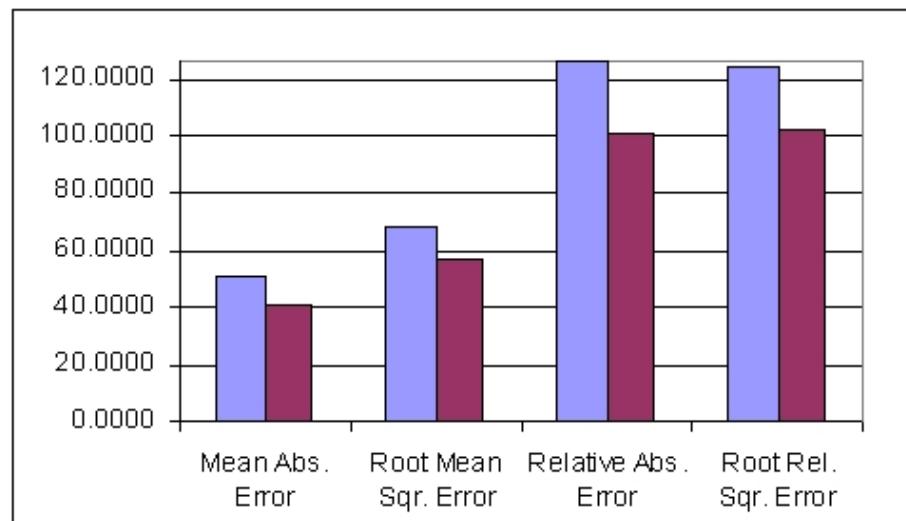
*Fig. 5. Maximum precision, recall and F-measure values over all the algorithms for the prediction of "Smoking" attribute, given only the social factors on a) Normal group, b) Pathologic group and c) Risk group, and given only the physical activity factors on d) Normal group, e) Pathologic group and f) Risk group.*

Let us see another representative example. Figure 6 presents the results of the prediction of cholesterol level from social factors, a), b), and c), and from physical activity factors, d), e) and f), for each of the level groups, respectively. In this case, the prediction results are very similar for the relationship between social factors and cholesterol, and physical activity and cholesterol, in all level groups, so we can conclude that the strength of the relationships is similar, too. However, it varies among level groups. In the Normal group, the mean absolute prediction error is about 24, being about 50 and 40 in Pathologic and Normal groups, respectively, concluding that it is easier to predict cholesterol, from both social factors and physical activity as training, for people in the Normal group. This fact denotes a strong relationship between the training factors and the cholesterol level in the latter group.

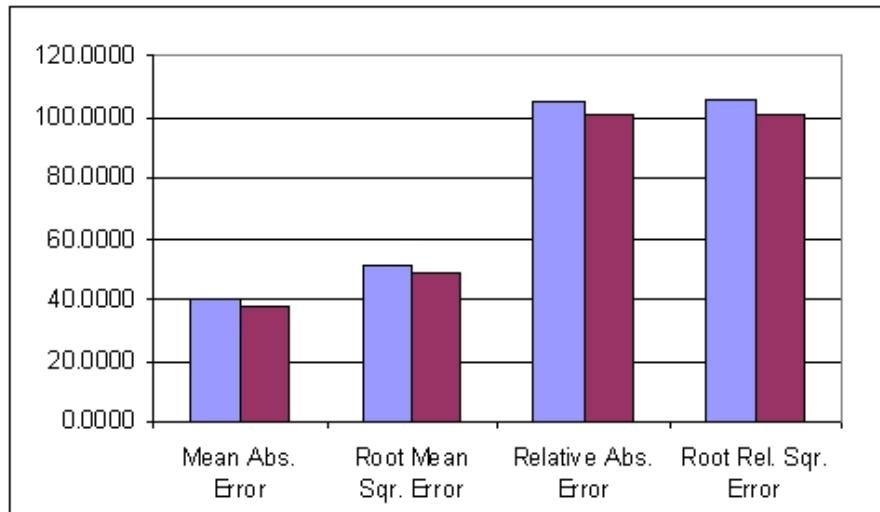
Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis



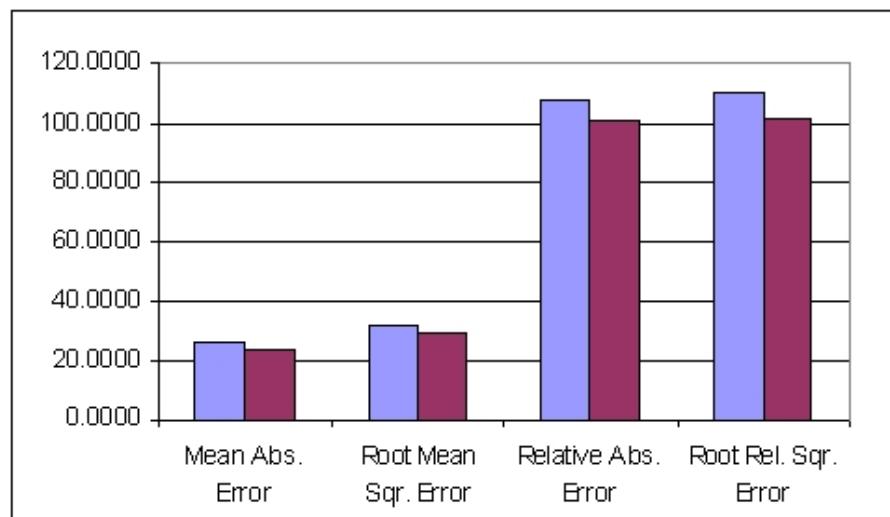
a)



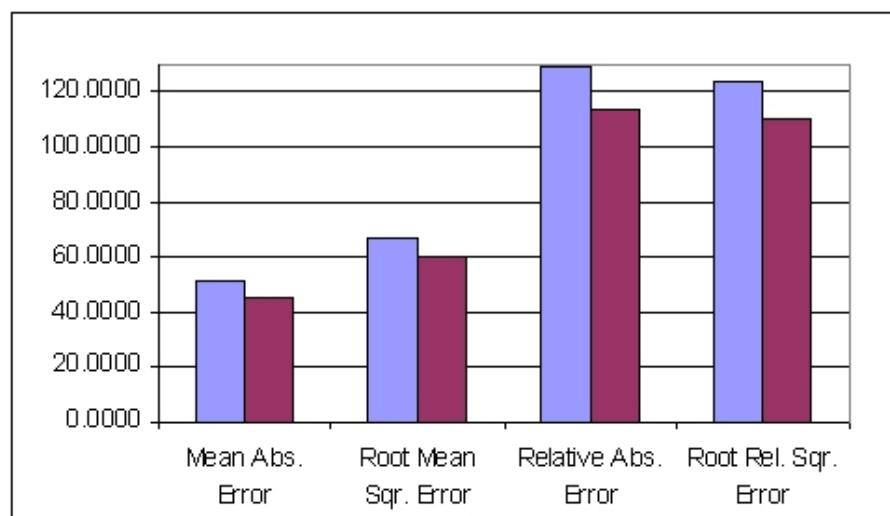
b)



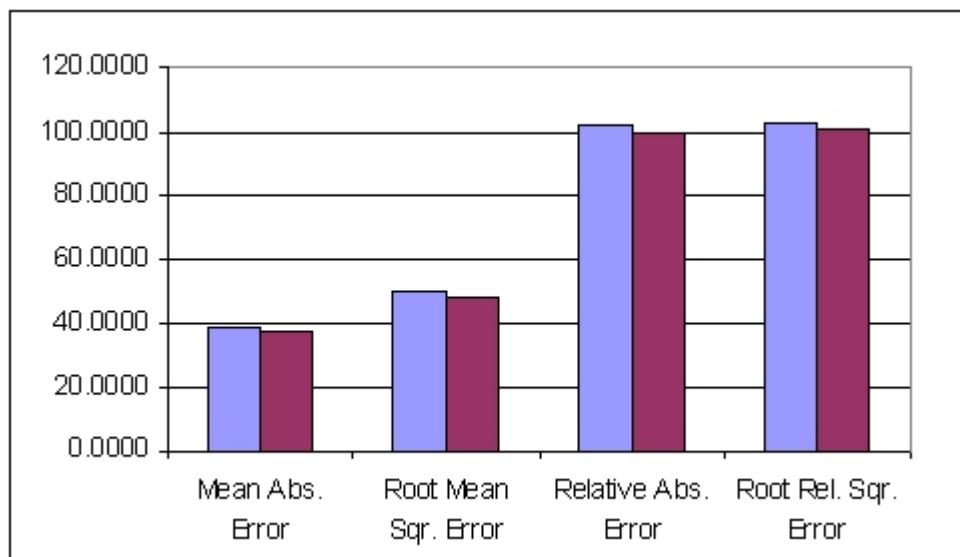
c)



d)



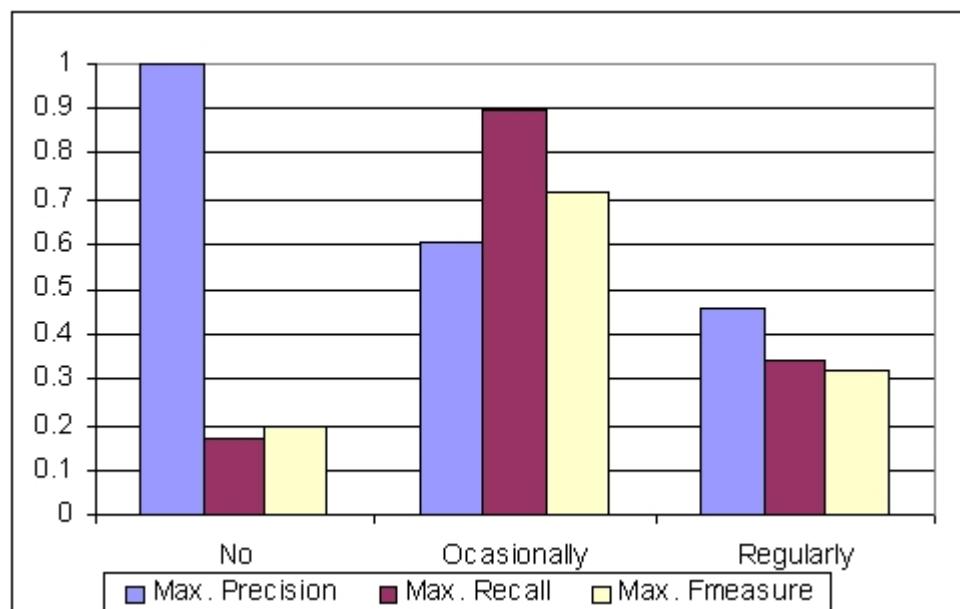
e)



f)

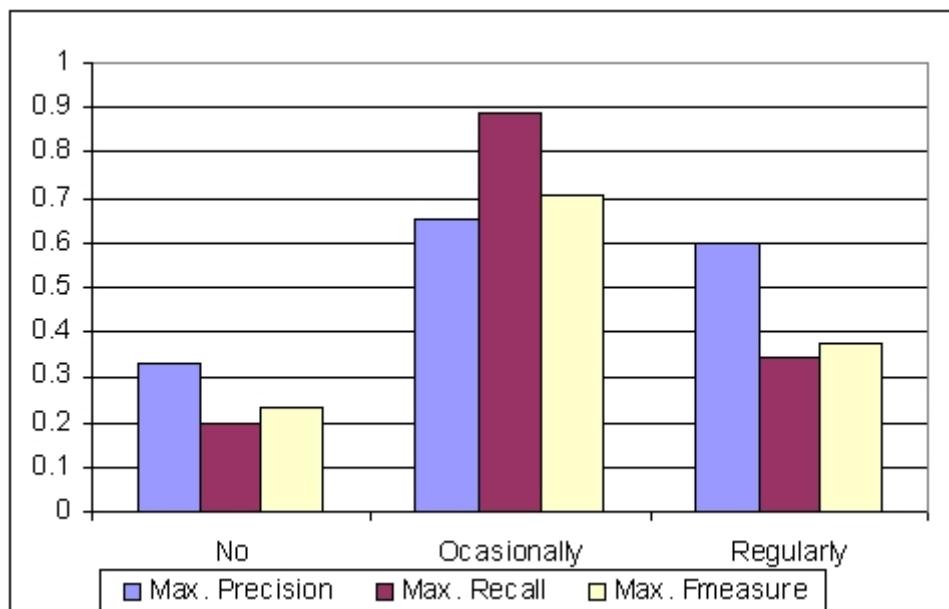
*Fig. 6. Average and maximum mean absolute error, root mean square error, relative absolute error and root relative square error values over all the algorithms for the prediction of cholesterol level attribute, given only social factors on a) Normal group, b) Pathologic group and c) Risk group, and given only the physical activity factors on d) Normal group, e) Pathologic group and f) Risk group.*

Finally, Figure 7 shows the results for the prediction of alcohol attribute values, separately from social factors and physical activity factors as training, for each level group, analogous to above.

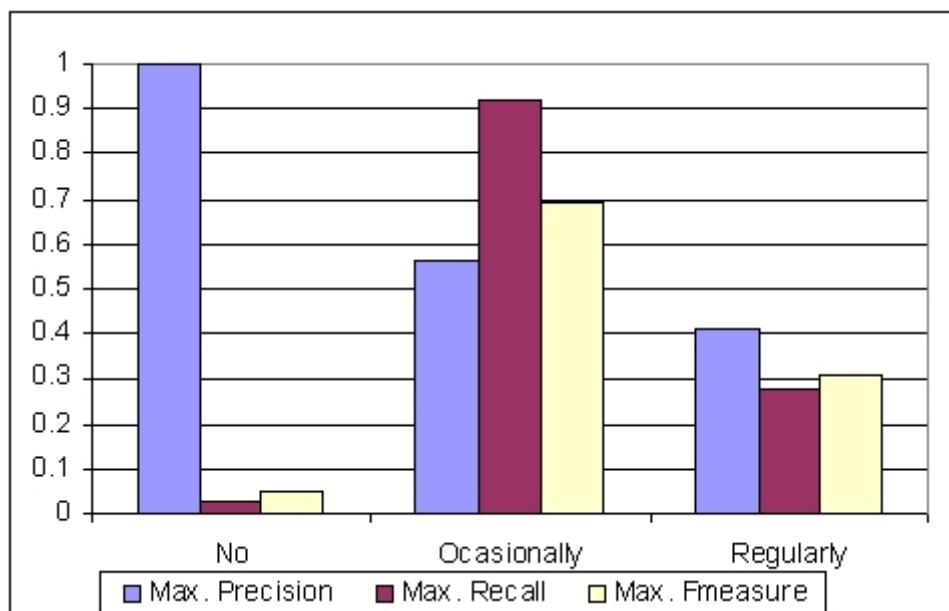


a)

Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis

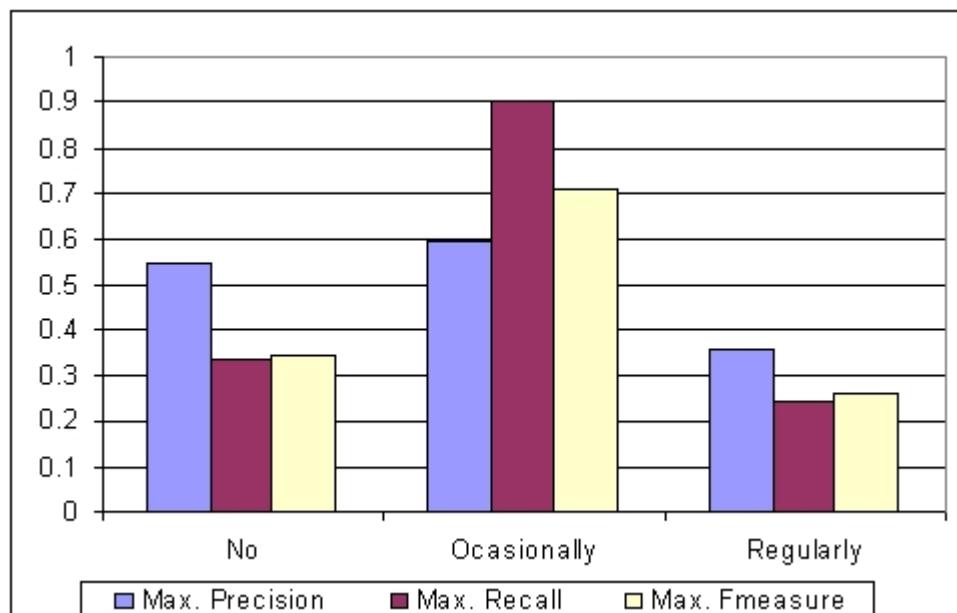


b)

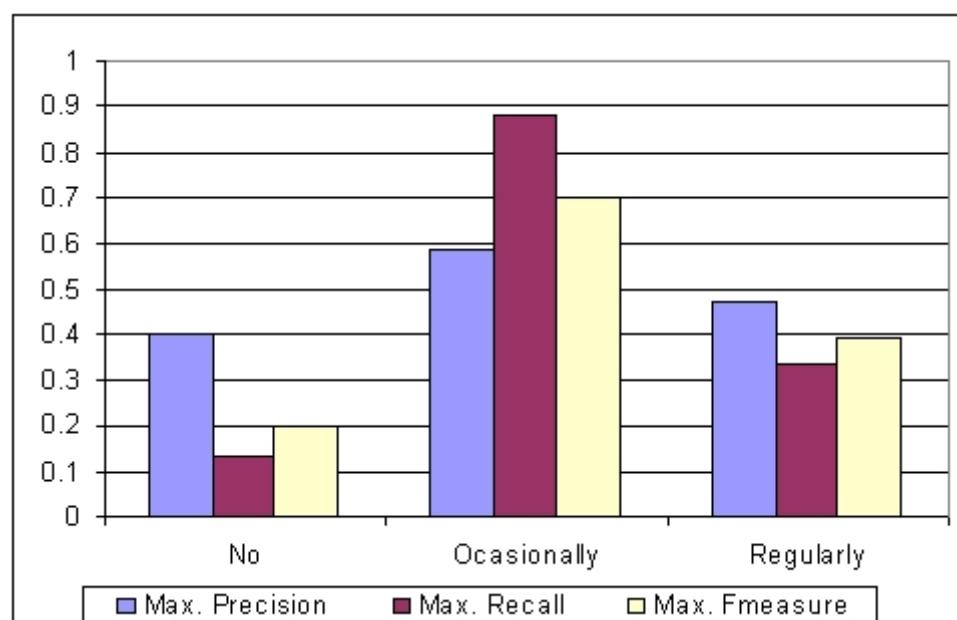


c)

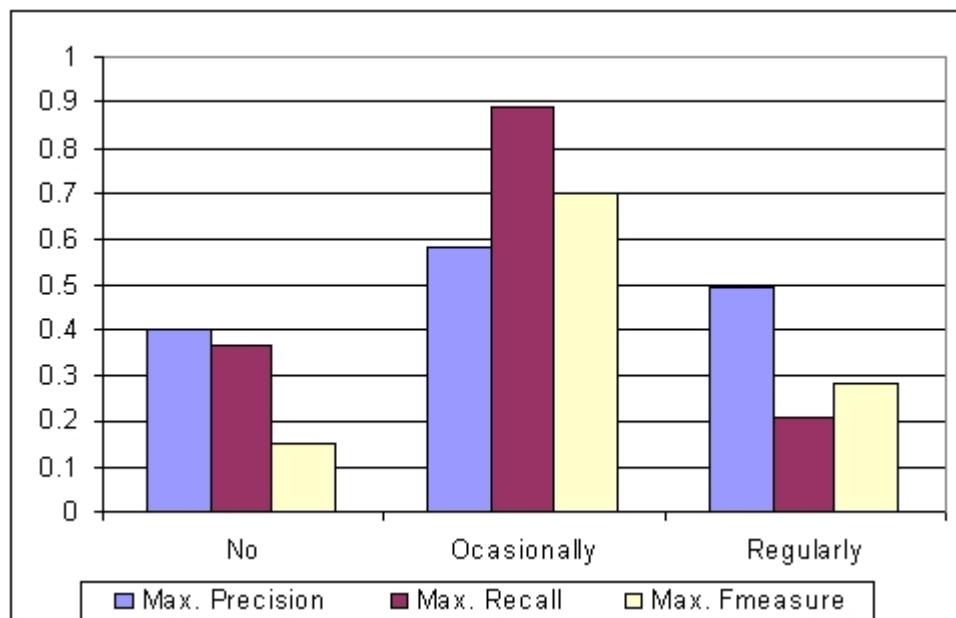
Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis



d)



e)



f)

*Fig. 7. Maximum precision, recall and F-measure values over all the algorithms for the prediction of the "Alcohol" attribute, given only the social factors on a) Normal group, b) Pathologic group and c) Risk group, and given only the physical activity factors on d) Normal group, e) Pathologic group and f) Risk group.*

The results in Figure 7 show that there is a clear relationship in all levels of groups between the training factors and the people who drink alcohol occasionally. People who drink regularly are more difficult to detect and predict from the training factors, resulting in a light relationship that is a little bit stronger in the Pathologic group. The same can be said about people who never drink alcohol in relation to physical activity factors. However, the prediction precision is significantly increased for social factors in Normal and Risk groups. People who never drink are accurately identified from their social factors in the latter groups, what denotes a significant relationship between the involved attributes.

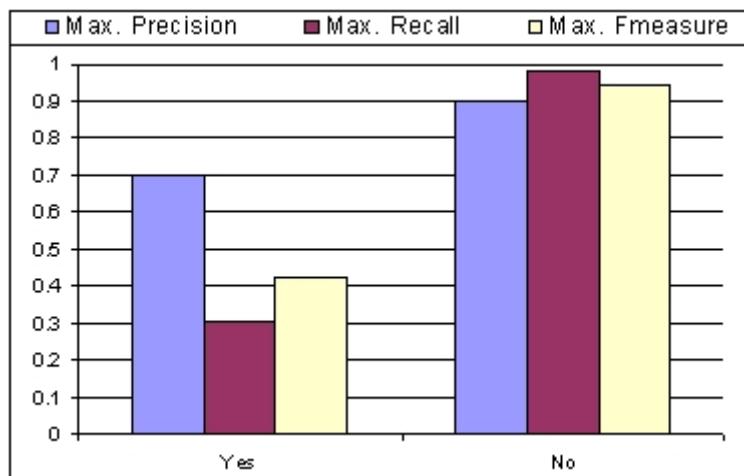
The training features groups are taken with all the attributes at once. From the medical point of view it is also interesting to separate these attributes and to try subsets of them. So, it was tried to predict the value of the physical activity in the job attribute given all possible combinations of social factors attributes, for example. The results show that, for Normal and Risk groups, the "Education" feature alone obtains much better prediction results than any other combination of social factors attributes. In the Pathologic group it is similar, but the difference is not so high as in the other groups, being "Age + Education" the best combination.

## 5.2 Predicting Future Disorders

The main objective of the next experiments is to test the prediction accuracy of the algorithms. The Entry collection is not the only one used but also the Control collection is considered. First of all, the patients who have a control record in the Control collection, after ten years from their entry in the study, were selected. Then using their Entry attributes it was

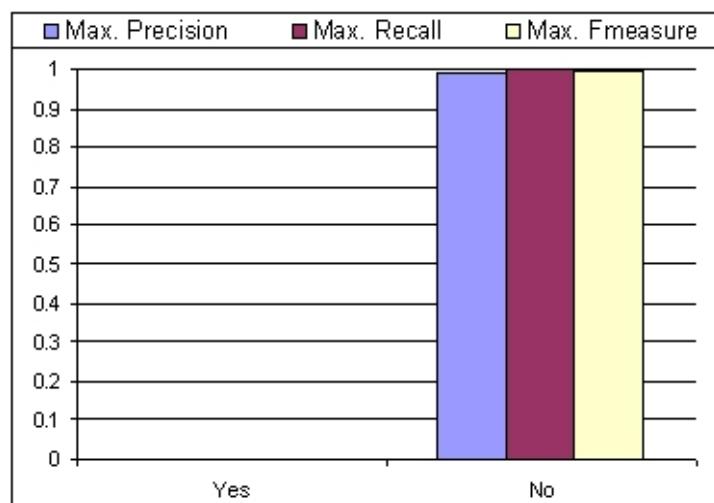
Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis

tried to predict whether they will have some disorders in ten years. These disorders correspond to systolic-diastolic hypertension, systolic hypertension, diastolic hypertension, hypercholesterolemia and hypertriglyceridemie. The possible values for these disorder attributes are true or false. The same has been done for twenty years records. The results show that the multilayer perceptron was the best algorithm, reaching values near 85 % of precision and 65 % of recall in the detection of all the disorders. The risk of future hypertension in the Risk group is 0 for many men, while some patients in this group were hypertensive since the beginning of the study. From the medical point of view it is more interesting to carry out the experiments only on the Normal group. The same process has been done separately for this group, for ten and twenty years. The results for the different mentioned disorders are presented in Figure 8, a) to e), respectively, for ten years prediction and f) to j), respectively, for twenty years prediction. For each disorder, the maximum values over all the different algorithms results are presented. In this case, results show that there is not one best algorithm. Depending on the disorder to predict and also on certain categories, one algorithm fits better than others (the maximum values presented correspond to different algorithms), so it will be interesting to use all the algorithms and make decisions based on the results from all of them. As a comment, we pointed that the prediction accuracy is much higher than when entries of the three levels of groups are considered all together, confirming the early interest in the Normal group.

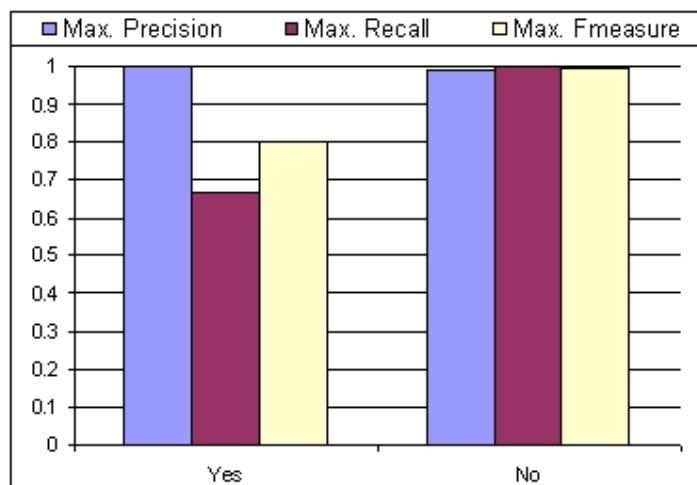


a)

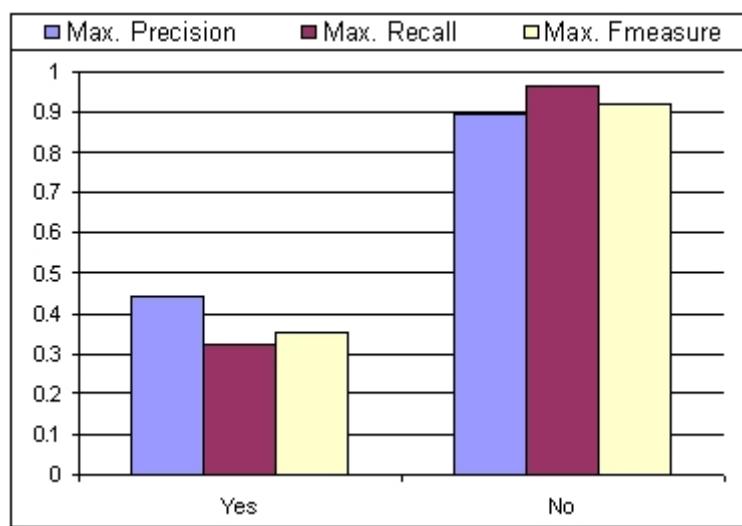
Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis



b)

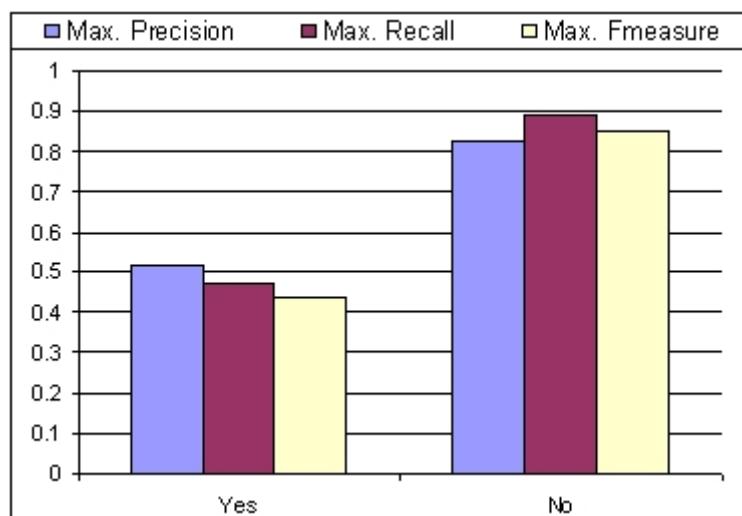


c)

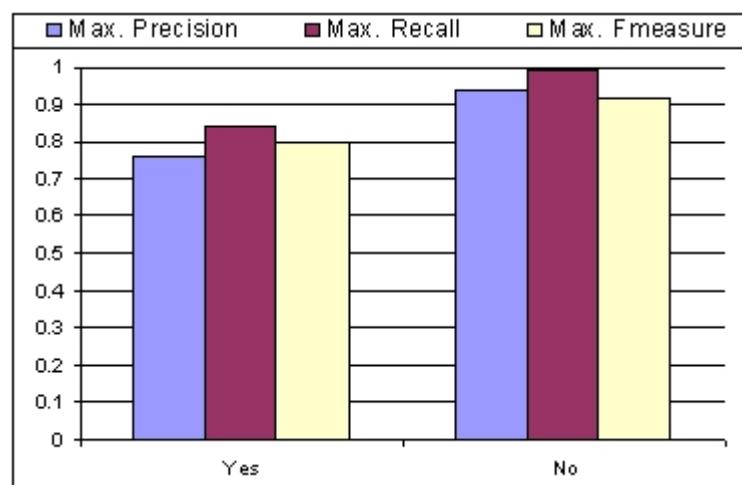


d)

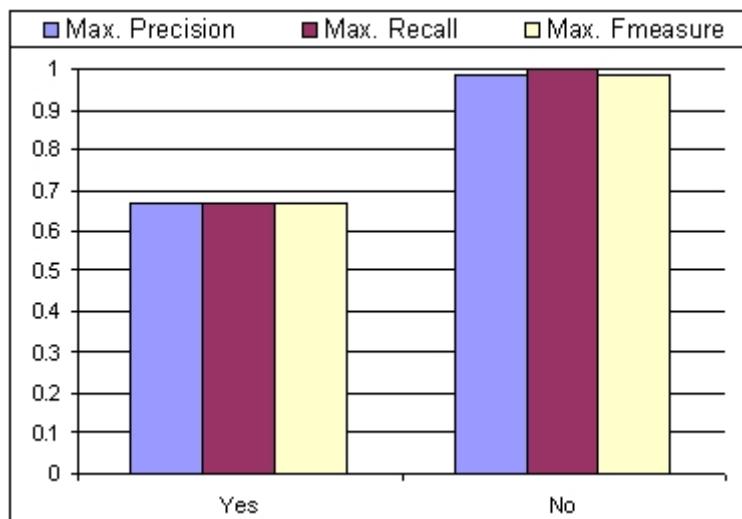
Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis



e)

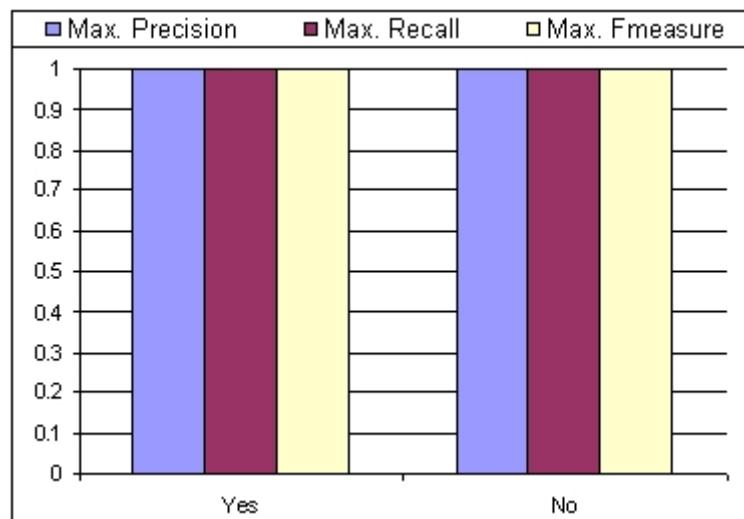


f)

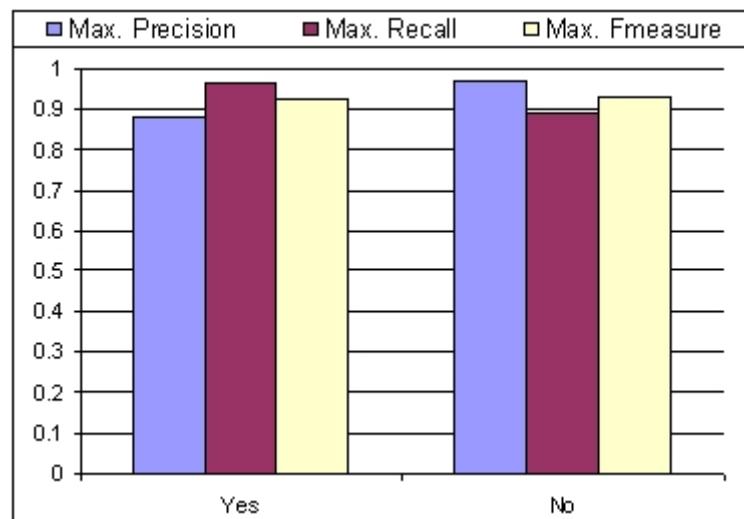


g)

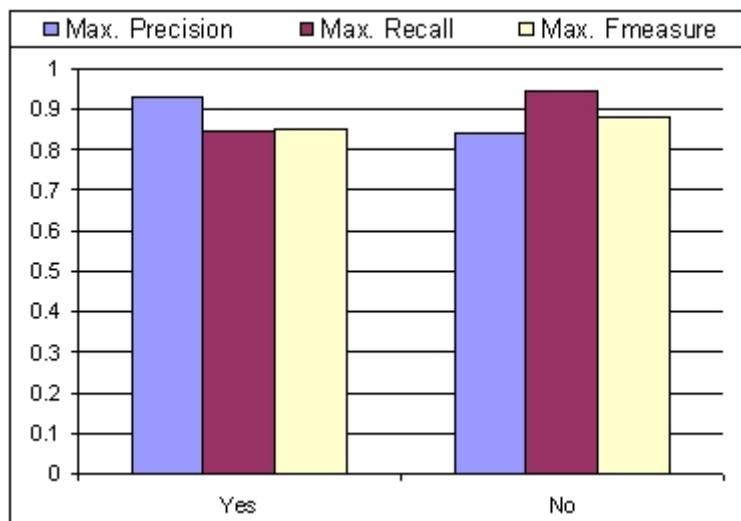
Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis



h)



i)



j)

*Fig. 8. Maximum precision, recall and F-measure values for the prediction of a) systolic-diastolic hypertension, b) systolic hypertension, c) diastolic hypertension, d) hypercholesterolemia and e) hypertriglyceridemie in ten years, and f) systolic-diastolic hypertension, g) systolic hypertension, h) diastolic hypertension, i) hypercholesterolemia and j) hypertriglyceridemie in twenty years.*

The values of Figure 8 show that it is more accurate to predict disorders in twenty years than to predict them in ten years, specifically the prediction of the presence of the disorders, which is accurately inferred in twenty years but very poorly predicted in ten years in all the disorders but diastolic hypertension. The non-presence of the disorders is equally well-predicted for both ten and twenty years. Among all the disorders the best detected is diastolic hypertension, obtaining prediction values near to 100 % accuracy for the presence and non-presence of the disorder. The worst predicted disorder is systolic hypertension, with the presence of the disorder non-detectable at all in ten years.

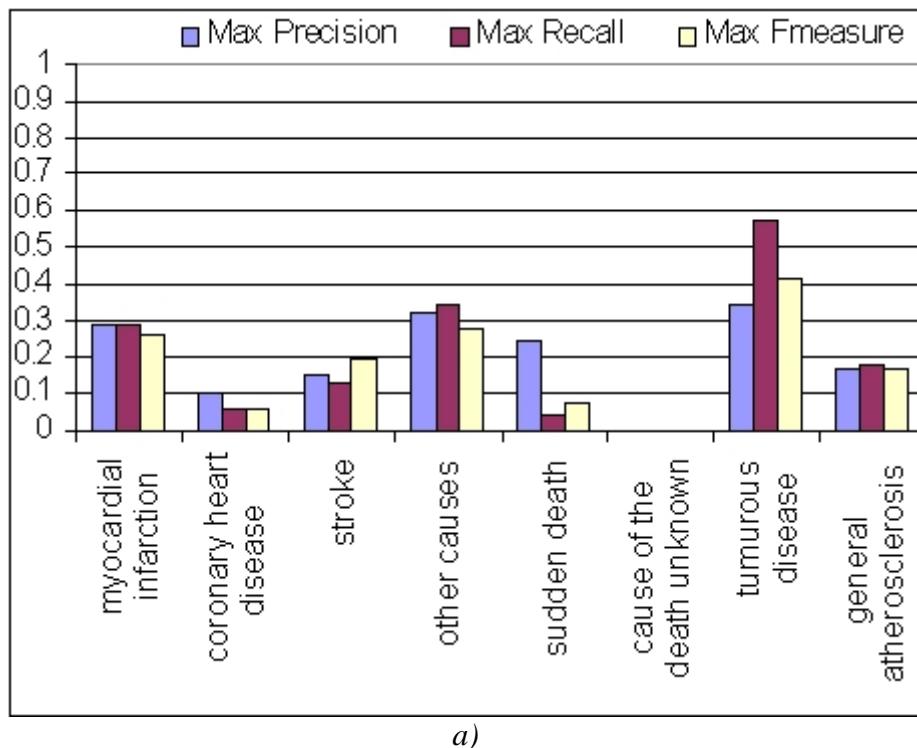
It was also tried to predict some other diseases, like angina pectoris, myocardial infarction, cerebro-vascular accident and so on, but there is a small number of observations with these features, so the results are not relevant.

### 5.3 Predicting Death Cause

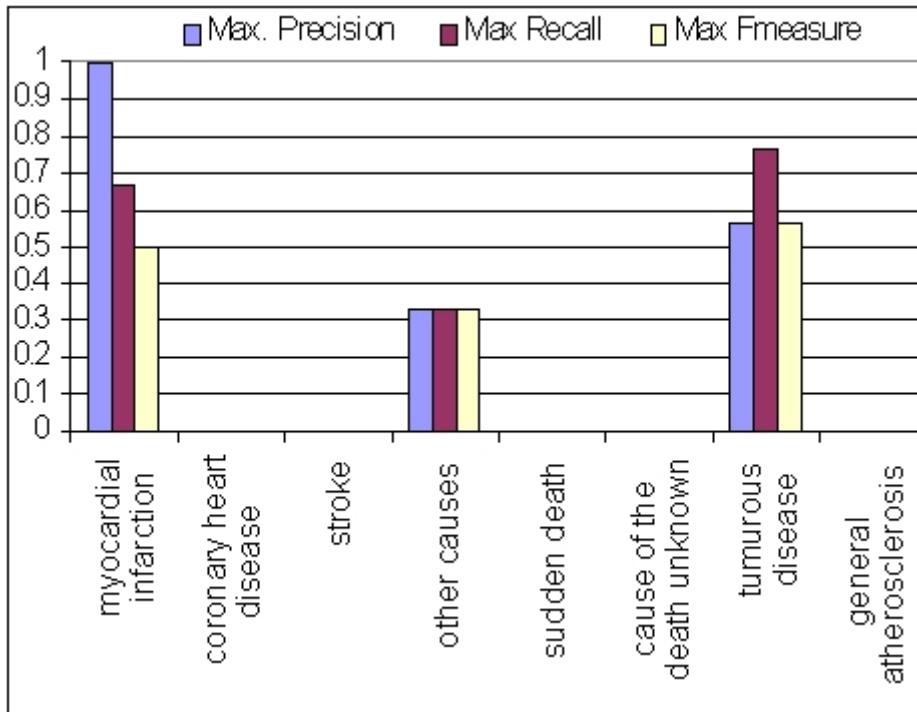
This experiment is analogous to the last one, but now what is tried is the prediction of the cause of death rather than diseases or disorders. Thus, the Death collection is used. The algorithms were trained with the data in the Entry collection for the patients of the Death collection. The experiments were carried out for the three levels of groups separately and for all the entries of all the groups together. The results are presented in Figure 9. In the Normal group, Figure 9b), the best predicted causes were tumour disease and other causes. In the Risk group, Figure 9d), the best prediction was for other causes but also for myocardial infarction and coronary heart disease, that were not predicted at all in the Normal group. In the Pathologic group, Figure 9c), the best predicted causes were tumour disease and myocardial infarction, but stroke and general atherosclerosis could be poorly predicted, too, obtaining much lower results for these latter causes in the other groups. In general, Figure 9a), the

Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis

results of prediction of death cause are very poor, concluding that data from the Entry collection has not enough information to predict death and/or also maybe more observations are needed. But, what is sufficient for it?

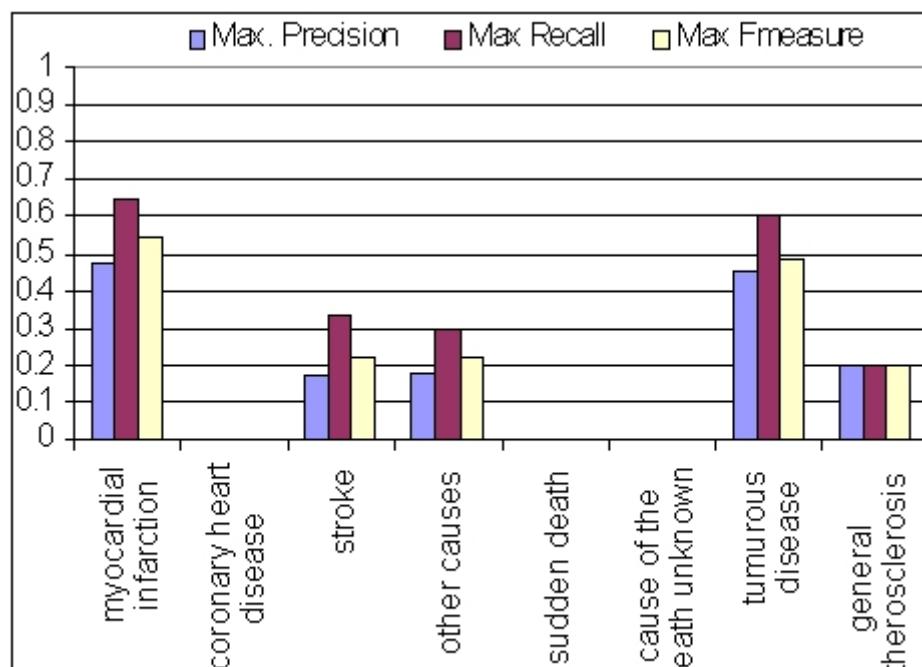


a)

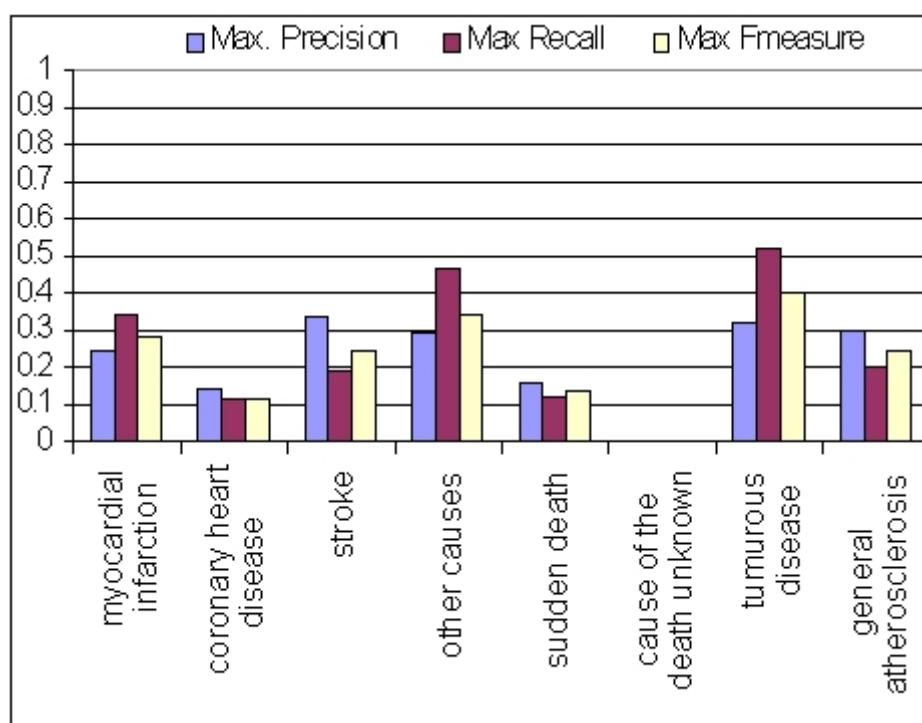


b)

Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis



c)



d)

Fig. 9. Maximum values of precision, recall and F-measure, in the prediction of death causes, for a) all the level groups as one, b) Normal group, c) Pathologic group and d) Risk group.

## 6. Conclusions

Machine learning algorithms belonging to a variety of paradigms have been applied to knowledge discovery on medical data in two different ways: firstly, the methods have been used in order to predict the value of one attribute of the patient database, given a subset of other attributes as training features, proposing the maximum accuracy among all the algorithms as a measure of the strength of the relationship between those training features and the target attribute. This measure has been proven useful also for comparing the relationships between attributes in different groups of patients.

Secondly, the learning techniques have been applied to the prediction of future disorders. The results show that some methods predict some disorders better than others. Then, it is interesting to use all the algorithms at a time and consider the result confidence based upon the known tendency of each method. All the tested methods perform better for twenty years prediction than for ten years predictions, reaching excellent results for some of the disorders that make the methods suitable for decision support. The machine learning algorithms have been also used in the prediction of death cause, obtaining poor results in this case, maybe due to the small amount of information (entries) of this type in the dataset.

It would be interesting for the future to finely tune the parameters of the algorithms and to test more techniques. It is also intended to integrate all methods with the degree of significance and usefulness discovered in this work in order to build an expert system, and the derivation of rules understandable by humans from the results of the system will be also researched.

## Acknowledgement

Research was partially supported by the Research Plan of ICS AS CR AV0Z10300504 and the “María Bueno” Research Visits Plan of the Spanish Council for Scientific Research together with the support of the Industrial Automation Institute.

## Reference

- [1] Mitchell, T.: Machine Learning. McGraw Hill, 1997.
- [2] Lavrać, N.: Selected Techniques for Data Mining in Medicine. Artificial Intelligence in Medicine, vol. 16 (1), pp. 3-23, 1999.
- [3] Aseervatham, S. and Osmani A.: Mining Short Sequential Patterns for Hepatitis Type Detection. ECML/PKDD Discovery Challenge, 2005.
- [4] Aubrecht, P., Kejkula, M., Kremen, P., Novakova, L., Rauch, J., Simunek, M., Stepankova, O.: Mining in Hepatitis Data by LISp-Miner and SumatraTT. ECML/PKDD Discovery Challenge, 2005.
- [5] Pizzi, L.C., Ribeiro, M.X., Vieira, M.T.P.: Analysis of Hepatitis Dataset using Multirelational Association Rules. ECML/PKDD Discovery Challenge, 2005.
- [6] Durand, N., Soulet, A.: Emerging Overlapping Clusters for Characterizing the Stage of Liver Fibrosis. ECML/PKDD Discovery Challenge, 2005.
- [7] Durand, N., Cleuziou, G., Soulet, A.: Discovery of Overlapping Clusters to Detect Atherosclerosis Risk Factors. ECML/PKDD Discovery Challenge, 2004.

Machine Learning Methods for Knowledge Discovery in Medical Data on Atherosclerosis

- [8] Cios, K. J.: Medical data mining and Knowledge Discovery. Physica – Verlag, 2001.
- [9] Chen, H., Fuller, S. S., Friedman, C. and Hersh, W.: Medical Informatics: Knowledge Management and Data Mining in Biomedicine. Integrated Series in Information Systems (2), Springer Science and Business Media Inc., 2005.
- [10] Boudik F., Reissigova J., Hrach K., Tomeckova M., Bultas J., Anger Z., Aschermann M., Zvarova J.: Primary Prevention of Coronary Artery Disease Among Middle Aged Men in Prague: Twenty-year Follow-up Results. *Atherosclerosis*. 2006 Jan;184(1):86-93.
- [11] Tomeckova, M.: The Challenge on Atherosclerosis Data Viewed by the Experts. ECML/PKDD Discovery Challenge, 2004.
- [12] Rish, I.: An Empirical Study of the Naive Bayes Classifier. IJCAI-01 Workshop on Empirical Methods in AI, 2001.
- [13] Haykin, S.: Neural Networks: A comprehensive Foundation (2nd edition). Pearson Education, 1998.
- [14] Scholkopf, B., Smola, A. J., Mtiller, K.-R., Burges, C. J. C., and Vapnik, V.: Support Vector Methods in Learning and Feature Extraction. In Down, T., Frean, M., and Gallagher, M., editors. Proceedings of the Ninth Australian Congress on Neural Networks, Brisbane, Australia. University of Queensland, 1998.
- [15] Teknomo, K.: K-Nearest Neighbors Tutorial. <http://people.revoledu.com/kardi/tutorialKNN>, 2004.
- [16] Quinlan, J. R.: C4.5: Programs for Machine Learning. Morgan Kauffman, 1993.
- [17] Compton, P., Edwards, G., Kang, B., Malor, R., Menzies, T., Preston, P., Srinivasan, A. and Sammut, S.: Ripple Down Rules: Possibilities and Limitations. Boose, J.H. & Gaines, B.R., Ed. Proceedings of the Sixth AAAI Knowledge Acquisition for Knowledge-Based Systems Workshop. pp.6-1-6-20. Calgary, Canada, University of Calgary, 1991.
- [18] Van Rijsbergen, C. J.: Information Retrieval. Butterworths, London, 1979.
- [19] Witten, I. H. and Frank, E.: Data Mining: Practical Machine Learning Tools and Techniques. 2nd Edition, Morgan Kaufmann, San Francisco, 2005.

# Métodos de Aprendizaje Automático para el Descubrimiento de Conocimiento en Datos Médicos sobre Arterosclerosis

José Ignacio Serrano<sup>1</sup>, Marie Tomečková<sup>2</sup>, Jana Zvárová<sup>2</sup>

1. Instituto de Automática Industrial, CSIC, Madrid, Spain,

2. Department of Medical Informatics, Institute of Computer Science AS CR, Prague, Czech Republic

Los algoritmos de aprendizaje automático son métodos que dado un conjunto de ejemplos de entrenamiento infieren un modelo de las categorías en las que se agrupan los datos, de tal forma que se pueda asignar a nuevos ejemplos una o más categorías de manera automática mediante analogía de patrones en dicho modelo. Los datos del estudio presentado son muy adecuados para este tipo de análisis. Muchos algoritmos de aprendizaje automático pertenecientes a todos los paradigmas han sido aplicados al descubrimiento de conocimiento en datos biomédicos. Los algoritmos utilizados en este trabajo pertenecen al paradigma del aprendizaje supervisado, tratando de cubrir la mayoría de las clases de algoritmos pertenecientes a este paradigma. Dos tipos de experimentos han sido realizados. Los del primer tipo tratan de validar los algoritmos en el descubrimiento de asociaciones entre atributos. El segundo tipo está orientado a la validación de la capacidad predictiva de disfunciones futuras. Los datos usados para estos experimentos han sido extraídos de un estudio preventivo, llevado a cabo durante 20 años, de los factores de riesgo de arterosclerosis en hombres de mediana edad. Este estudio se denomina STULONG (LONGitudinal STUdy). Los resultados obtenidos muestran que algunos algoritmos predicen ciertas enfermedades mejor que otros, por lo que es interesante usar todos ellos al mismo tiempo y ponderar sus resultados individuales en base a la eficacia conocida de cada uno de ellos con respecto a la enfermedad objetivo. Se ha probado también la capacidad predictiva de los algoritmos sobre las causas de muerte, obteniendo resultados pobres debido quizás a la escasa información presente en los datos sobre estas causas.

**Palabras clave:** descubrimiento de conocimiento, aprendizaje automático supervisado, minería de datos biomédicos, factores de riesgo de arterosclerosis

## 1. Introducción

Los algoritmos de aprendizaje automático son métodos que dado un conjunto de ejemplos de entrenamiento infieren un modelo de las categorías en las que se agrupan los datos, de tal forma que se pueda asignar a nuevos ejemplos una o más categorías de manera automática mediante analogía de patrones en dicho modelo. Estas técnicas han sido aplicadas con éxito a una gran variedad de problemas y datos en tareas de predicción. El objetivo principal de este trabajo es investigar y descubrir cómo aplicar los algoritmos de aprendizaje supervisado para descubrir relaciones entre atributos y para realizar predicciones que puedan ser útiles a la toma de decisiones. Los datos biomédicos son un tipo especial de datos, ya que datos de

diferente naturaleza recogen toda la información. Además, este tipo de datos presenta ciertos problemas conocidos: información ausente y dispersa, ruido y datos temporales. Los algoritmos de aprendizaje automático son muy apropiados para este tipo de datos [2]. Existen algunos trabajos sobre KDD que intentan tratar con información biomédica a gran escala. En [3], los autores tratan de detectar el tipo de hepatitis extrayendo patrones de secuencia corta a partir de características temporales. En [4], se generan reglas sencillas usando *4ft-miner* (es decir, tablas estadísticas de dos filas y dos columnas) para caracterizar las diferencias temporales existentes entre las hepatitis B y C. Los autores de [5] tratan de descubrir reglas de atributos binarios sencillos que sean capaces de predecir el nivel de fibrosis del hígado. El mismo objetivo se persigue en [6] pero utilizando patrones que son posteriormente agrupados y asignados a niveles de fibrosis. Esta misma técnica es también aplicada a la detección del riesgo de arterosclerosis en [7]. Otros ejemplos de minería de datos biomédicos se presentan en [8] y [9]. Para los experimentos del trabajo aquí presentado se ha utilizado la colección denominada STULONG (LONGitudinal STUdy) [10], resultado de un estudio durante 20 años de los factores de riesgo de arterosclerosis en hombres de mediana edad. El principal objetivo del trabajo aquí presentado es validar los algoritmos de aprendizaje automático como un método de descubrimiento de asociaciones considerando la eficacia de clasificación como una medida de importancia de las asociaciones extraídas. Además, se trata de comprobar la capacidad de predicción de enfermedades futuras de los algoritmos.

En la siguiente sección se presentan los detalles de la colección STULONG. En la sección 3, se describen los algoritmos de aprendizaje automático empleados. En la sección 4 se presentan las medidas de evaluación del comportamiento de los algoritmos y en la sección 5 se muestran los experimentos y sus resultados. Finalmente, algunas conclusiones y trabajo futuro se exponen en la sección 6.

## 2. Descripción del estudio y de la colección de datos

El corpus STULONG [10] [11] fue recopilado por el 2<sup>nd</sup> Departamento de Medicina Interna, y la 1<sup>st</sup> Facultad de Medicina y el Hospital General Facultativo de Praga, y transformado a formato electrónico, así como analizado estadísticamente por el Centro Europeo de Informática Médica, Estadística y Epidemiología (EuroMISE) de la Universidad Charles y la Academia de Ciencias de la República Checa.

Las principales pretensiones del estudio fueron:

1. Identificar la presencia de factores de riesgo (RF) de arterosclerosis en una población generalmente considerada como la más afectada por posibles complicaciones de la enfermedad en hombres de mediana edad.
2. Seguir el desarrollo de estos factores de riesgo y su impacto en la salud de los sujetos examinados, especialmente respecto a las enfermedades cardiovasculares.
3. Estudiar el impacto de la intervención compleja de los factores de riesgo en el desarrollo de enfermedades cardiovasculares y en la mortalidad por dichas causas.

Los hombres nacidos entre 1926 y 1937 que vivían en le distrito 2 de Praga fueron seleccionados para el estudio en 1975. Para el primer reconocimiento, 1419 de los 2370 sujetos invitados se prestaron al estudio. La invitación incluía una pequeña explicación de los

objetivos del mismo. A los sujetos que declinaron la participación se les envió otras dos invitaciones para reconocimientos posteriores.

Los factores de riesgo fueron definidos en términos de niveles de la siguiente manera:

- hipertensión – presión sanguínea BP  $\geq 160/95$  mm Hg u hombres bajo medicación hipertensiva,
- hipercolesterolemia - colesterol  $\geq 260\text{mg\%}$  (6,7 mmol/l),
- hipertrigliceridemia – triglicéridos  $\geq 200\text{mg\%}$  (2,2 mmol/l),
- fumador:  $\geq 15$  cig./día actualmente o la misma cantidad hasta un año previo al estudio (los fumadores de puros o pipas no fueron considerados como fumadores),
- sobrepeso: índice de Brocka  $> 115\%$  (índice Brocka: altura en cm menos 100 = 100 %),
- historia médica familiar no favorable: muerte del padre o la madre por enfermedad arterial coronaria o ataque cardíaco anterior a los 65 años de edad.

De acuerdo a la presencia de los factores de riesgos anteriores, estado de salud general y resultados de ECG, los sujetos fueron divididos en los siguientes grupos:

- NG** = grupo de sujetos sin factores de riesgo, sin manifestación de enfermedades arteriales u otras enfermedades severas que hicieran imposible su observación durante los 10 años siguientes, y sin cambios ECG.
- RG** = grupo de sujetos con al menos un factor de riesgo, sin manifestación de enfermedades arteriales u otras enfermedades severas que hicieran imposible su observación durante los 10 años siguientes, y sin cambios ECG.
- PG** = grupo de sujetos con una enfermedad cardiovascular manifiesta u otro tipo de disfunción severa que hace imposible su observación en los siguientes 10 años. Este grupo patológico (PG) incluye también a sujetos con diabetes tratada con fármacos o insulina y a sujetos con ECG patológico, de acuerdo al código ECG de Minnesota.

Las observaciones a largo plazo de los pacientes se realizaron de acuerdo a los grupos descritos anteriormente:

- El grupo de riesgo **RG** fue dividido aleatoriamente en dos subgrupos designados como **RGI** (grupo de riesgo intervenido) y **RGC** (grupo de riesgo de control). Los pacientes en el grupo **RGI** fueron invitados al reconocimiento un mínimo de dos veces al año. Siguiendo la administración farmacológica fueron reconocidos cuando fue necesario. Los pacientes en el grupo **RGC**

recibieron un pequeño informe escrito que contenía sus resultados de laboratorio y su descripción ECG, además de una recomendación de presentar estos resultados a sus doctores. La intervención con respecto a estos resultados fue puesta en manos de sus doctores. En un primer estudio, no se encontró ninguna diferencia significativa en edad, factores socioeconómicos o factores de riesgo entre ambos grupos.

- El 10 % de los sujetos en el grupo **NG** fue examinado un mínimo de una vez al año. – (se les denomina **NGS**); En este grupo, análogamente al grupo de riesgo, la intervención fue iniciada en cuanto se identificó y se confirmó alguno de los factores

de riesgo (hiperlipidemia, hipertensión arterial, etc.). El resto de sujetos del grupo NG fueron invitados a realizar un control entre 10 y 12 años después.

- Los sujetos del grupo **PG** fueron excluidos de cualquier observación posterior.

La intervención fue un problema clave en el estudio y se basó en la influencia no farmacológica. Se trataron de modificar y optimizar los factores de riesgo:

- Intervención No Farmacológica*: recomendaciones sobre estilo de vida, dieta, actividad física, hábito de fumar, disminución de peso, etc. Las recomendaciones fueron realizadas en cada observación centradas en un factor de riesgo específico en cada sujeto, exceptuando algunas instrucciones de carácter general.
- Intervención Farmacológica*: tratamiento de la hipertensión arterial y de la hiperlipoproteinemia. Fue muy leve en las etapas tempranas del estudio y más severa en los últimos años del mismo. La terapia farmacológica fue impuesta con respecto al conjunto general de factores de riesgo.

Después de todo el proceso de adquisición de los datos, cuatro conjuntos de los mismos fueron usados para el análisis:

1. El conjunto *ENTRY* contiene valores de 244 atributos obtenidos de la primera observación de cada sujeto. Estos atributos son códigos o medidas y transformaciones de medidas de diferentes variables (identificador, familia e historial personal, factores sociales – educación, actividad física, tabaco, hábitos de dieta, alcohol, medidas antropométricas – altura, peso, presión sanguínea, pulso, pruebas de laboratorio y código ECG).
2. El conjunto *CONTROL* contiene resultados de exámenes de control con 66 atributos. Estos atributos corresponden al identificador, cambios en los hábitos, historial personal, valores físicos y bioquímicos, y datos sobre hipertensión, hipercolesterolemia, hipertrigliceridemia y otras enfermedades coronarias y oncológicas. Este conjunto se compone de un total de 10,572 registros.
3. Información adicional sobre el estado de salud de 403 sujetos que declinaron la primera invitación fue recogida mediante cuestionarios por correo. Los valores de los 62 atributos recogidos en los cuestionarios se almacenaron en el conjunto *LETTER*.
4. El conjunto *DEATH* contiene información sobre las causas de la muerte de 389 pacientes, descritas mediante 5 atributos, además del identificador de los sujetos y la fecha de su muerte.

### 3. Descripción de las técnicas de aprendizaje automático empleadas

Todos los algoritmos empleados pertenecen al paradigma de aprendizaje supervisado, es decir, necesitan una etapa de aprendizaje para construir un modelo a partir de los datos de entrenamiento para después usar ese modelo en la predicción o inferencia de la categoría de ejemplos desconocidos. Se han empleado varios algoritmos tratando de representar a todas las clases de algoritmos dentro del paradigma. Cada uno de estos algoritmos se describe brevemente a continuación:

### 3.1 Naïve Bayes

Naïve Bayes [12] calcula, para cada par atributo-valor, por ejemplo (*educación, universitaria*), la probabilidad de pertenecer a cada categoría, dividiendo el número de ejemplos de cada categoría donde el par aparece entre el número total de ejemplos donde el par aparece. De esta forma, cada par tendrá asociado una probabilidad para cada posible categoría. Naïve Bayes está basado en la suposición de que cada par atributo-valor de un ejemplo es independiente del resto. Así, cuando un ejemplo nuevo se clasifica, la probabilidad asociada a cada categoría es la multiplicación de la probabilidad para la correspondiente categoría de cada uno de los pares que conforman el ejemplo. La categoría final asignada es la que más alta probabilidad asociada tiene.

### 3.2 Perceptrón Multicapa

El modelo de clasificación de la Red Neuronal Perceptrón Multicapa [12] está compuesto de un cierto número de capas de neuronas interconectadas entre sí. La arquitectura usada en este trabajo se muestra en la Figura 1.

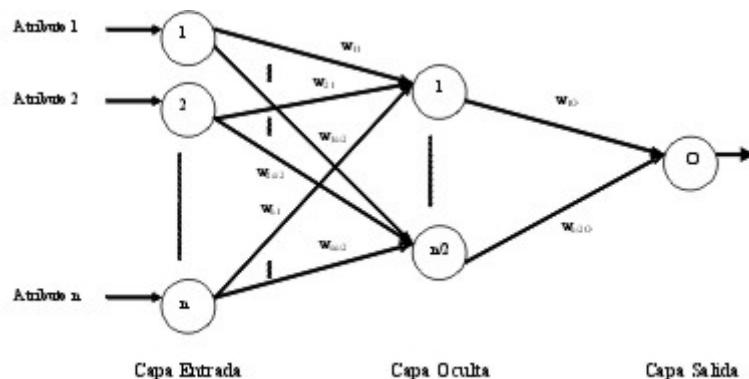


Figura 1. Arquitectura empleada de la Red Neuronal Perceptrón Multicapa.

Cada conexión tiene un peso asociado. La entrada a cada neurona es la suma ponderada, usando los pesos de asociación, de todos los valores entrantes. La salida de cada neurona es el resultado de aplicar una función. En este caso, se implementó una función sigmoide típica en todas las neuronas. La Figura 2 muestra la representación y expresión de la función.



Figura 2. Expresión y representación de la función sigmoide.

Así, cada valor de cada atributo de un ejemplo se introduce en la neurona correspondiente de la capa de entrada y los valores se propagan por la red hasta la capa de salida, donde el valor resultante de la neurona de salida corresponde a la categoría inferida.

La etapa de entrenamiento consiste en, dado un conjunto de pesos iniciales, introducir cada uno de los ejemplos de entrenamiento en el modelo y comparar el valor de salida con la categoría real esperada. Dependiendo del error de la clase inferida, el algoritmo *backpropagation* modifica los pesos desde la capa de salida a la de entrada para hacer que la salida inferida sea igual a la esperada. Este proceso se lleva a cabo un número determinado de épocas o iteraciones. En este caso, se emplean 500 iteraciones en la fase de entrenamiento. La cantidad en la que los pesos son modificados, llamada tasa de aprendizaje, es igual a 0.3 y el momento aplicado a los pesos durante su actualización es de 0.2. Si el algoritmo de aprendizaje no alcanza una buena aproximación a la salida esperada después de una iteración completa, se inicializa y provoca un decremento en la tasa de aprendizaje.

### 3.3 Máquinas de Vectores de Soporte (SVM)

Las Máquinas de Vectores de Soporte [14] intentan separar los ejemplos, basándose en su categoría, en el espacio de  $n$  dimensiones siendo  $n$  el número total de atributos o características, mediante hiperplanos de la forma  $\mathbf{w} + \mathbf{b}$ , tal que

$$\mathbf{x} \mathbf{w} + b \geq +1 \rightarrow \text{categoría} = \text{sí}$$

$$\mathbf{x} \mathbf{w} + b \geq -1 \rightarrow \text{categoría} = \text{no}$$

siendo  $\mathbf{x}$  el ejemplo representado como un vector de  $n$  componentes. Aquí,  $\mathbf{w}$  es el vector de soporte perpendicular al hiperplano, y corresponde a los ejemplos que se sitúan más allá o en los límites de la categoría a la que pertenecen (ver Figura 3).

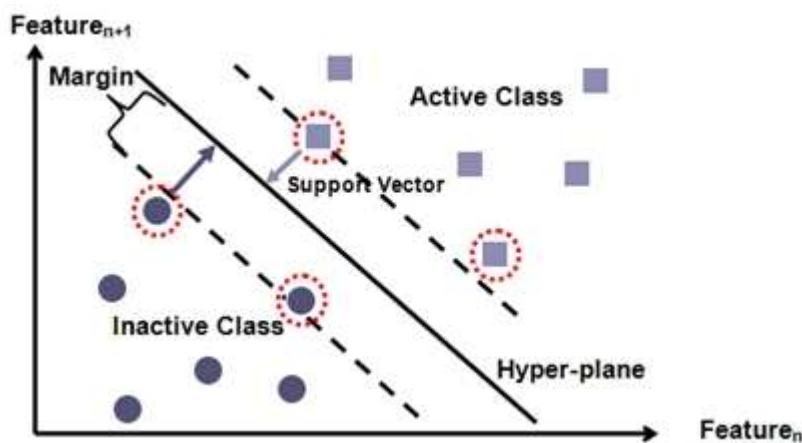


Figura 3. Esquema de los Vectores de Soporte.

Los vectores de soporte definen también, mediante su módulo, un margen unitario entre el hiperplano y los ejemplos positivos y negativos más cercanos (esa es la razón de los umbrales +1 y -1). Para cada categoría el algoritmo trata de encontrar  $\mathbf{w}$  maximizando el margen. Para clasificar un ejemplo nuevo simplemente se aplica la expresión anterior. Esta simple implementación del método es la que se emplea en los experimentos, aunque existe una gran abanico de variaciones mucho más sofisticadas.

### 3.4 K-Vecinos Más Cercanos (KNN)

KNN es un algoritmo basado en memoria [15], con la idea subyacente de que las experiencias pasadas pueden ayudar a resolver las presentes mediante analogía. Considera a cada ejemplo como un vector de  $n$  componentes, siendo nuevamente  $n$  el número de atributos o características. No necesita una etapa de aprendizaje. Para inferir la clase de un ejemplo desconocido hasta el momento, el algoritmo compara ese ejemplo con todos los ejemplos de entrenamiento o memoria calculando la distancia entre ellos. A continuación, la clase mayoritaria de entre los  $K$  ejemplos más similares al de entrada es la categoría inferida para el mismo. La medida de distancia empleada es la distancia Euclídea entre dos vectores. Sin embargo, existen más posibilidades recogidas en la literatura.

### 3.5 Árboles de Decisión ID3 y C4.5

El modelo producido por este algoritmo es un árbol [16], donde cada nodo corresponde a un atributo y cada arco del nodo corresponde a un posible valor del atributo nodo.

El algoritmo de aprendizaje construye el árbol a partir de los datos de entrenamiento. La selección del atributo que formará un nodo en cada momento es llevada a cabo mediante el cálculo de la entropía de los datos después la selección. Esto es, para cada atributo se calcula la entropía de los datos restantes agrupados por los posibles valores del atributo evaluado. El atributo cuyos valores produzcan una entropía menor es el seleccionado para formar el siguiente nodo. El proceso continua hasta que no hay más atributos que seleccionar o bien hasta que el número de ejemplos agrupados bajo un nodo es menor que un umbral. En este último caso, se forma un nodo hoja correspondiente a la categoría mayoritaria de los nodos agrupados bajo ese nodo. En la Figura 4 se muestra un ejemplo sencillo:

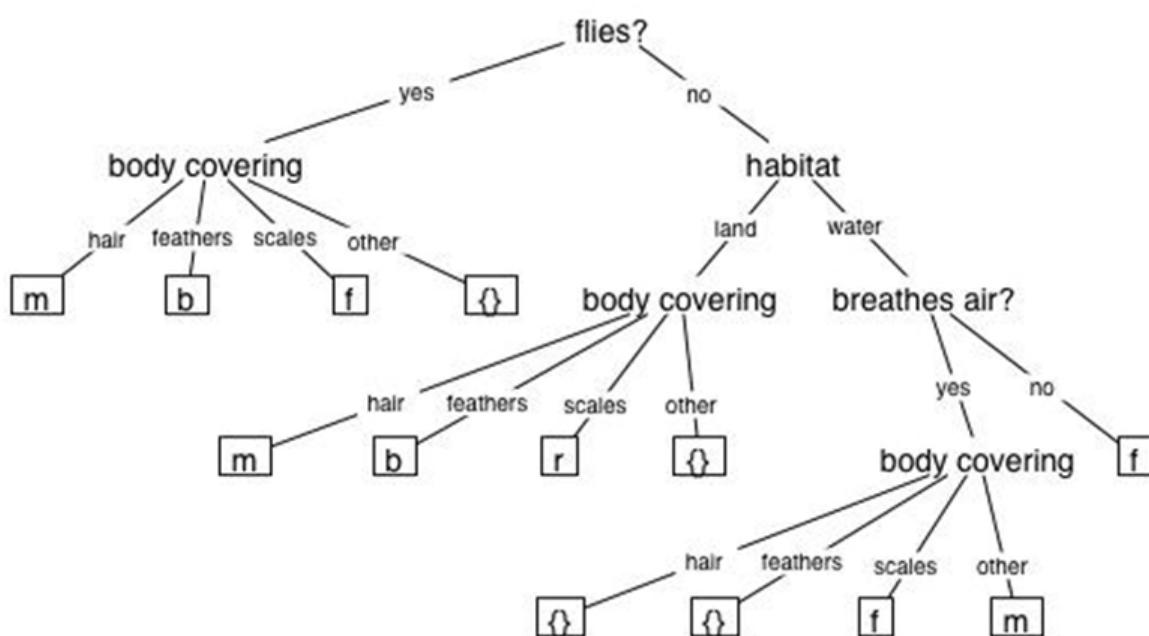


Figura 4. Ejemplo de árbol de decisión.

En el ejemplo se identifican 4 atributos: *vuela, recubrimiento del cuerpo, hábitat y respira aire*, y cuatro posibles categorías: *m, b, f y r*. En este caso, el primer atributo es *vuela* porque es el que produce la división de los datos con entropía mínima en ese nivel, y análogamente con el resto. Para clasificar un ejemplo nuevo sólo hay que seguir el árbol de arriba abajo y la hoja final es la categoría inferida. Los caminos desde la raíz hasta los nodos hoja se pueden ver como reglas, donde el antecedente está formado por la intersección de los pares atributo-valor de los caminos.

C4.5 es una ampliación de ID3 que permite el uso de atributos numéricos continuos, tiene en cuenta los valores ausentes y realiza un proceso de poda inteligente del árbol para reducir su tamaño y permitir así tratar con un gran número de ejemplos. El árbol J48 usado en los experimentos de este trabajo es una implementación de C4.5.

### 3.6 Extracción de Reglas Ridor

Ridor abrevia a RIpple-DOWN Rule [17]. Este algoritmo genera una regla por defecto que se ajuste a la mayoría de los ejemplos de entrenamiento y luego busca excepciones con la menor tasa de error al clasificar los propios ejemplos de entrenamiento. A continuación genera las excepciones a las excepciones con menos error, de manera recursiva. Así, lleva a cabo una expansión de excepciones en forma de árbol donde la raíz está formada por la regla por defecto. Las excepciones son un conjunto de reglas que predicen las clases que no contempla la regla por defecto. IREP una implementación de Ridor y es el algoritmo empleado para encontrar excepciones. Éste construye las reglas añadiendo un término al antecedente en cada iteración de tal forma que el error se minimice. Los términos del antecedente son de la forma (*atributo {=, ≠, ≤, ≥} valor*).

## 4. Evaluación

Los procesos y medidas de evaluación son los mismos para todos los experimentos: dada la colección de datos, una parte de la misma es considerada como conjunto de entrenamiento y el resto como conjunto de test. Así, los modelos aprenden del conjunto de entrenamiento y tratan de inferir las categorías de los ejemplos del conjunto de test. Puesto que las categorías de éstos últimos son conocidas se pueden validar las inferencias de los modelos. Así pues, esta validación se realiza para cada categoría mediante tres medidas típicas: precisión, cobertura y medida-F [18]. La precisión es el porcentaje de predicciones de una categoría que son correctas. La Ecuación 1 muestra la expresión analítica de la precisión.

$$\text{Precisión}(\text{categoría}_i) = \frac{n^{\circ} \text{ de predicciones correctas de categoría}_i}{n^{\circ} \text{ total de predicciones de categoría}_i} \quad (1)$$

La cobertura es el porcentaje de todos los ejemplos de test pertenecientes a una categoría y que son correctamente inferidos. Su expresión analítica se muestra en la Ecuación 2.

$$\text{Cobertura}(\text{categoría}_i) = \frac{n^{\circ} \text{ de predicciones correctas de categoría}_i}{n^{\circ} \text{ total de ejemplos de categoría}_i} \quad (2)$$

La medida-F es una combinación de las medidas anteriores. Representa, en cierto modo, la intersección entre los ejemplos implicados en la precisión y la cobertura, normalizada mediante la suma de ambas. La Ecuación 3 presenta su expresión analítica.

$$\text{Medida-F} = \frac{2 * \text{Precisión} * \text{Cobertura}}{\text{Precisión} + \text{Cobertura}} \quad (3)$$

Así pues, estas tres medidas se calculan para cada categoría del conjunto de test. Como se ha comentado antes, dada la colección es necesario dividirla en conjuntos de entrenamiento y de test. Un proceso común de evaluación es la validación cruzada. La colección se divide en  $n$  partes iguales. A continuación, cada combinación de  $n-1$  partes se emplea como conjunto de entrenamiento y la parte restante como test, de tal forma que el algoritmo se ejecuta  $n$  veces y las medidas finales son la media de las  $n$  ejecuciones. Para todos los experimentos descritos a continuación el valor de  $n$  es igual a 3, de tal forma que el entrenamiento siempre es un 66% del conjunto total, ejecutando cada algoritmo 3 veces. Habitualmente, el valor de  $n$  es mayor que 3 (típicamente igual a 10), pero en este caso se tienen pocos ejemplos en algunas de las categorías y un valor mayor de  $n$  podría generar conjuntos de test sin representación en las categorías mencionadas, lo que no es en absoluto deseable.

## 5. Experimentos

Se han llevado a cabo dos tipos de experimentos. En el primero de ellos se tratan de descubrir asociaciones entre atributos considerando las medidas de la eficacia de la clasificación como un indicador de la importancia de las asociaciones. El resto de experimentos están encaminados a la validación de los algoritmos en la predicción de enfermedades futuras.

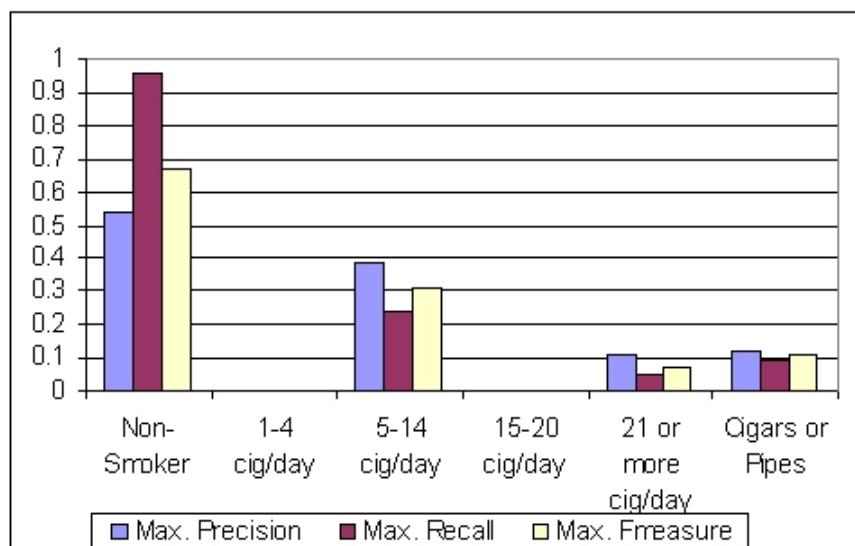
Es necesario indicar que las observaciones en los datos que presentan valores ausentes no fueron eliminadas ni sustituidas, ya que las implementaciones de los algoritmos empleados son capaces de tratar con dichos valores. Dichas implementaciones son las incluidas en el entorno de descubrimiento de conocimiento WEKA [19], usadas con los parámetros por defecto en los experimentos descritos a continuación.

### 5.1 Encontrando respuestas

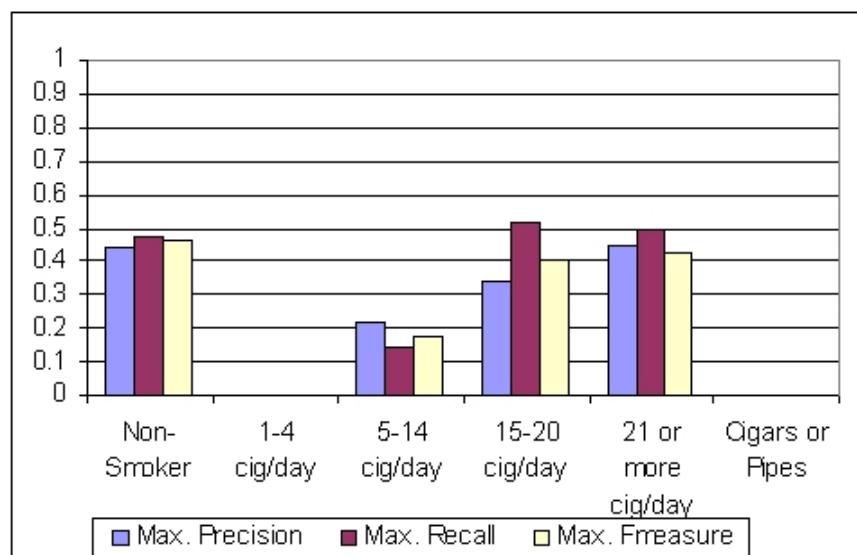
Los primeros experimentos están relacionados con las cuestiones analíticas propuestas por el Discovery Challenge de la conferencia ECML/PKDD 2004, y más concretamente con las relacionadas con el conjunto ENTRY. Estas tareas consisten en el descubrimiento de relaciones en tres grupos diferentes de sujetos: grupo Normal (NG), grupo de Riesgo (RG) y grupo Patológico (PG). Estos grupos corresponden a los niveles de riesgo de arterosclerosis descritos en la Sección 2 y serán denominados como grupos de nivel de ahora en adelante. Las relaciones objetivo son las que asocian los atributos referentes a factores sociales con el resto, los referentes a actividad física con el resto y así sucesivamente. De esta manera, los algoritmos de aprendizaje automático son aplicados de manera independiente a los datos correspondientes a cada grupo de nivel, intentando predecir el valor de cada uno de los atributos como categorías. Así, por ejemplo, dados los cuatro atributos correspondientes a factores sociales como atributos de entrenamiento, los algoritmos se ejecutan de manera independiente para predecir el valor de cada uno de los cuatro atributos de actividad física, y

análogamente con el resto de atributos. Para cada relación, se calculan los valores máximos de entre los resultados de todos los algoritmos. Estos valores permitirán la comparación entre grupos de nivel. Si la eficacia de clasificación es buena se puede decir que existe una relación fuerte, con un grado igual al valor de dicha eficacia, entre los atributos empleados para el entrenamiento y el atributo cuyos valores se tratan de inferir. Dicho grado permitirá comparar las relaciones en los distintos grupos de nivel.

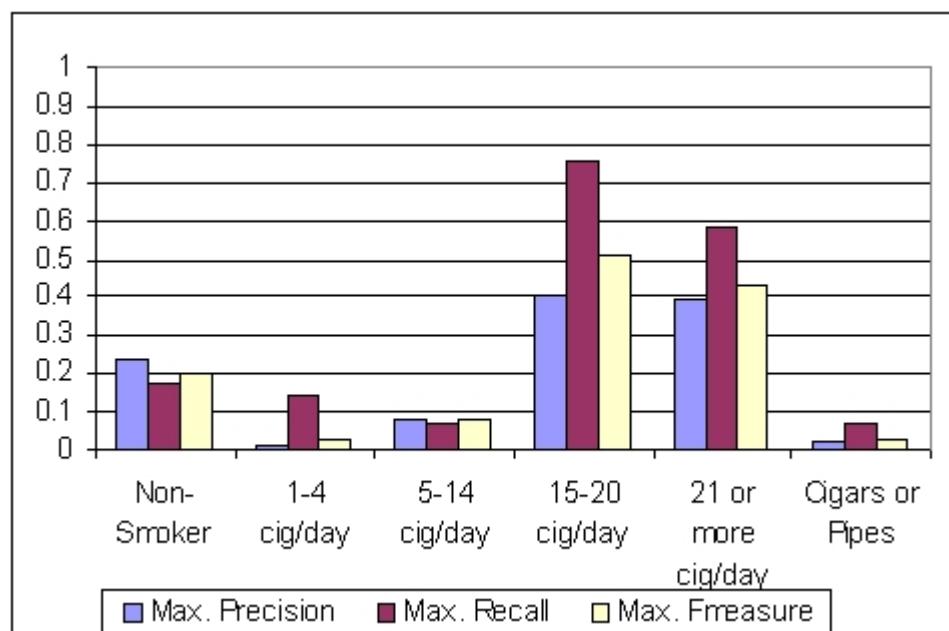
Debido a las limitaciones de espacio, sólo se presentan algunos de los resultados más representativos. En la Figura 5 se presenta, para cada uno de los grupos de nivel (a) Normal, b) Patológico y c) Riesgo, respectivamente, la precisión, cobertura y medida-F máximas para la predicción del atributo “Fumar” dados los factores sociales, y análogamente dados los factores físicos de d) a e). Como se puede observar, en el grupo Normal, los mejores resultados de predicción se obtienen para los sujetos no fumadores, ya sea a partir de los factores sociales o de la actividad física, siendo el resto de valores no significativos. Es perceptible también que la relación entre factores sociales y “Fumar” es ligeramente más fuerte que entre actividad física y “Fumar”, ya que todas las medidas son algo mejores en el primer caso. Tanto en los grupos Patológico como de Riesgo, la relación entre los atributos de entrenamiento y la categoría de no fumadores es más fuerte con la actividad física, siendo especialmente fuerte en el grupo Patológico. En estos grupos, los sujetos que fuman 15 o más cigarrillos al día son más precisamente detectados que en el grupo Normal, aunque no así los no fumadores.



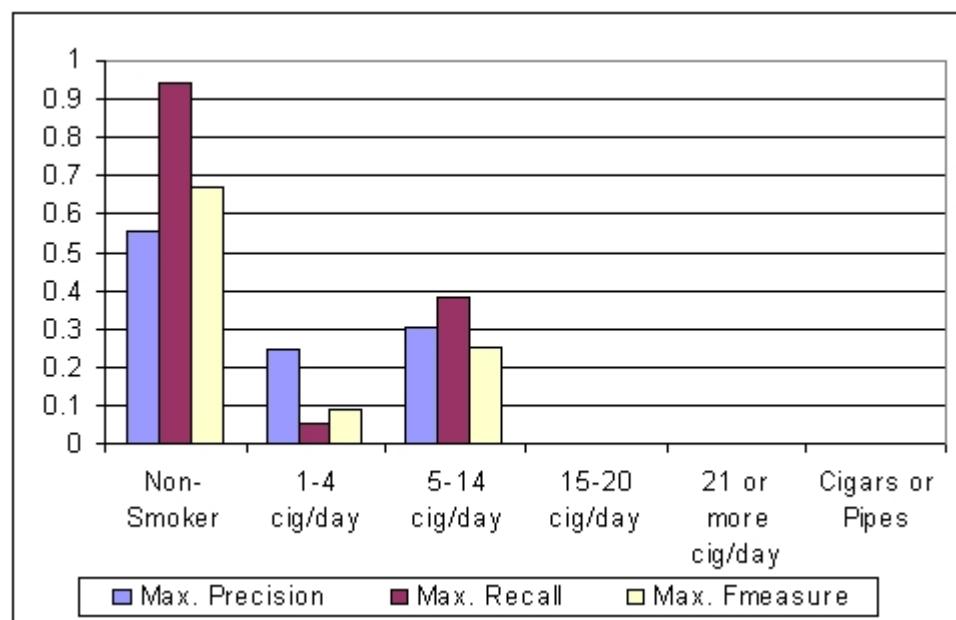
a)



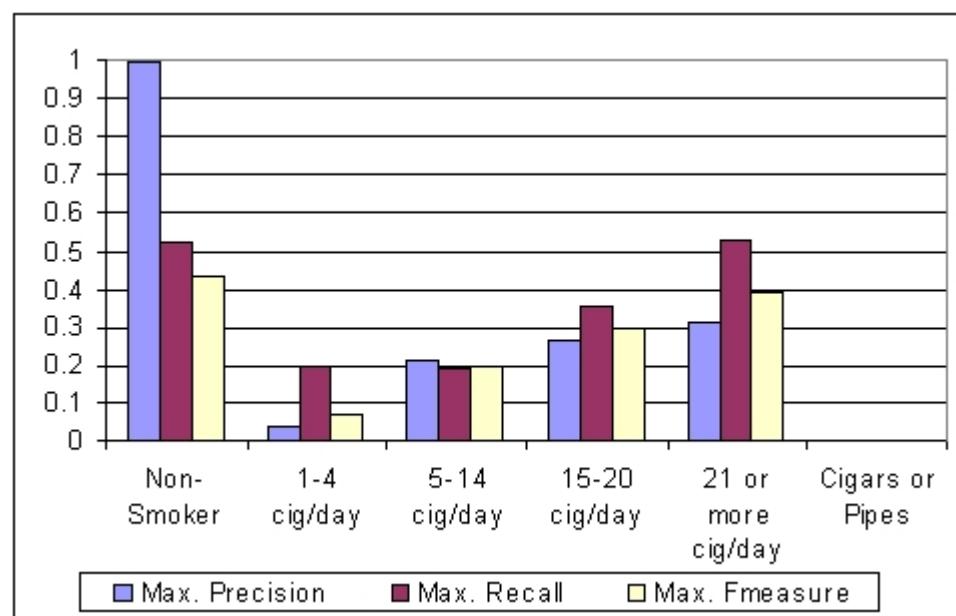
b)



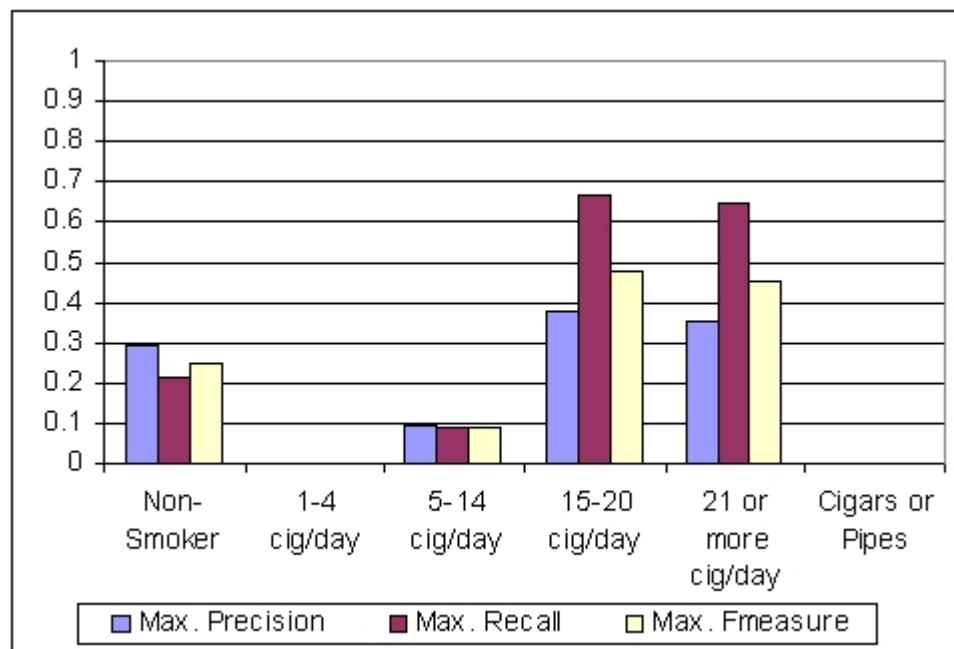
c)



d)

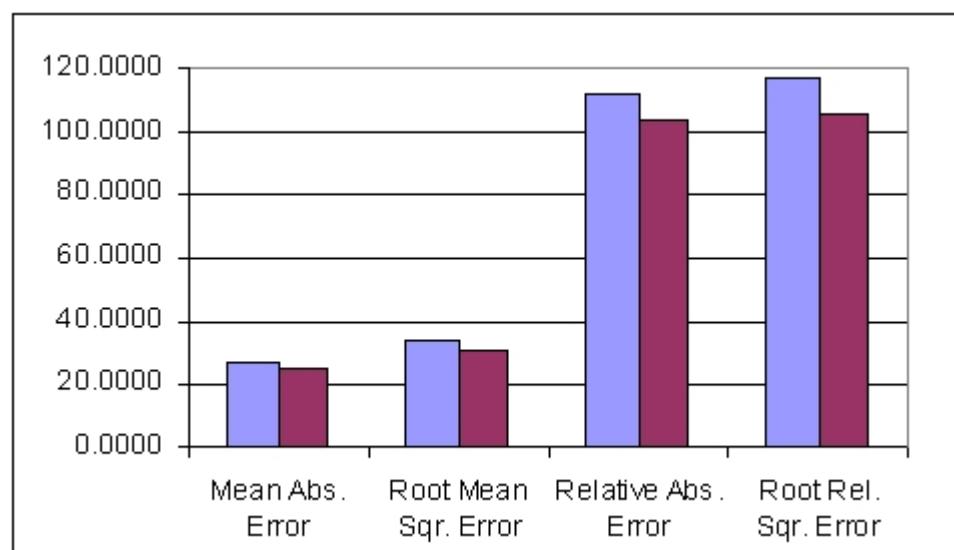


e)

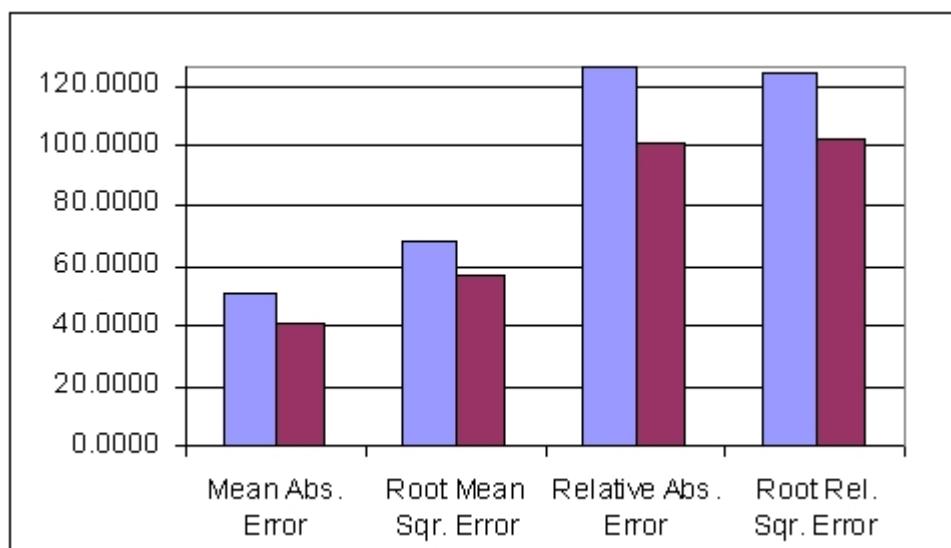


f)

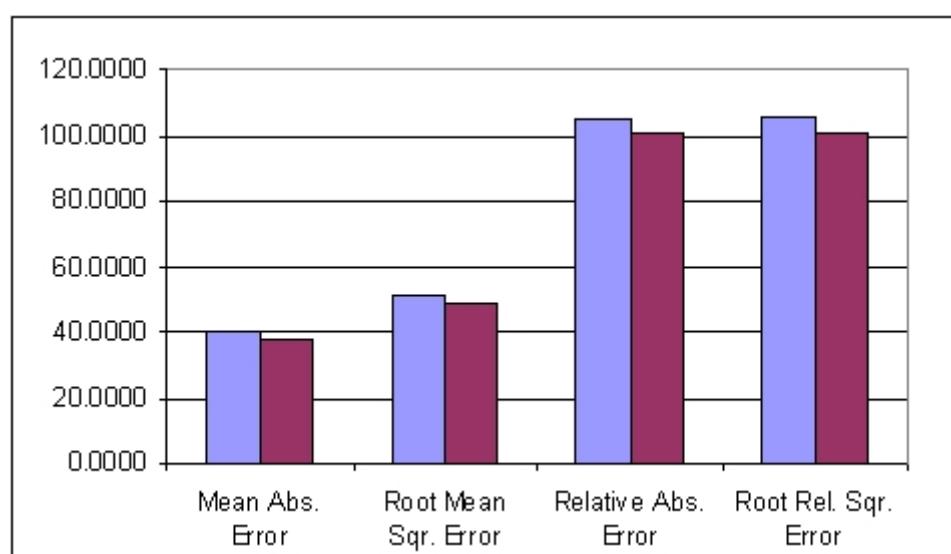
Figura 5. Precisión, cobertura y medida-F máximas de entre todos los algoritmos para la predicción de "Fumar", dados los factores sociales en a) grupo Normal, b) grupo Patológico y c) grupo de Riesgo, y dada la actividad física en d) grupo Normal, e) grupo Patológico y f) grupo de Riesgo.



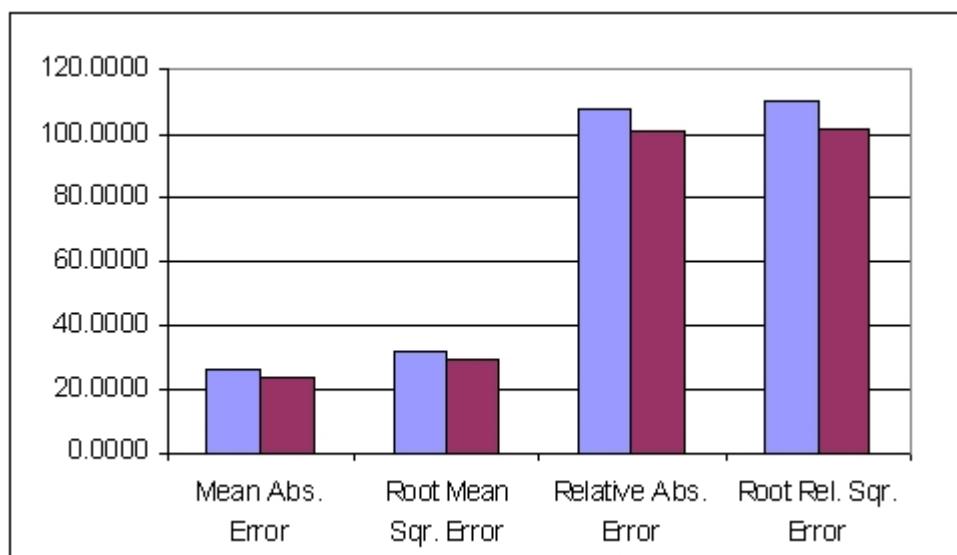
a)



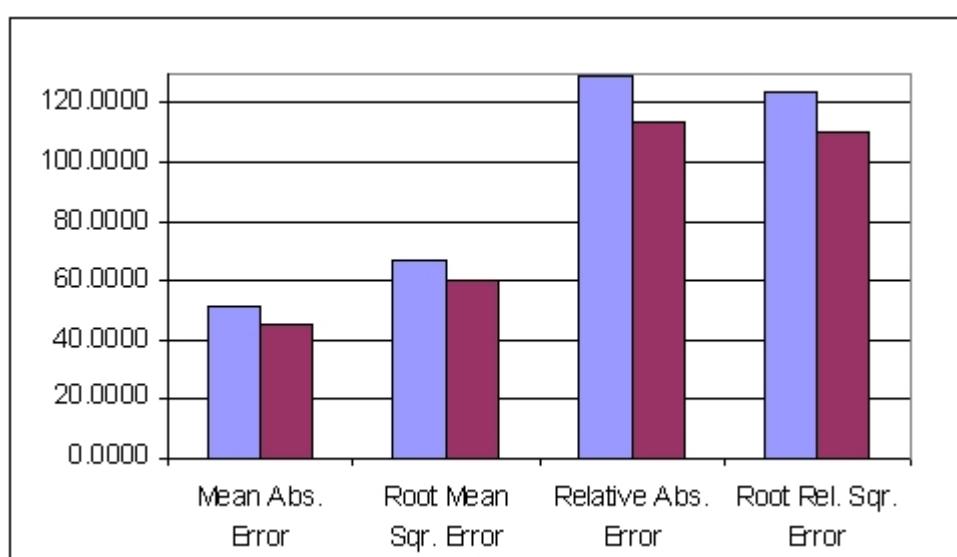
b)



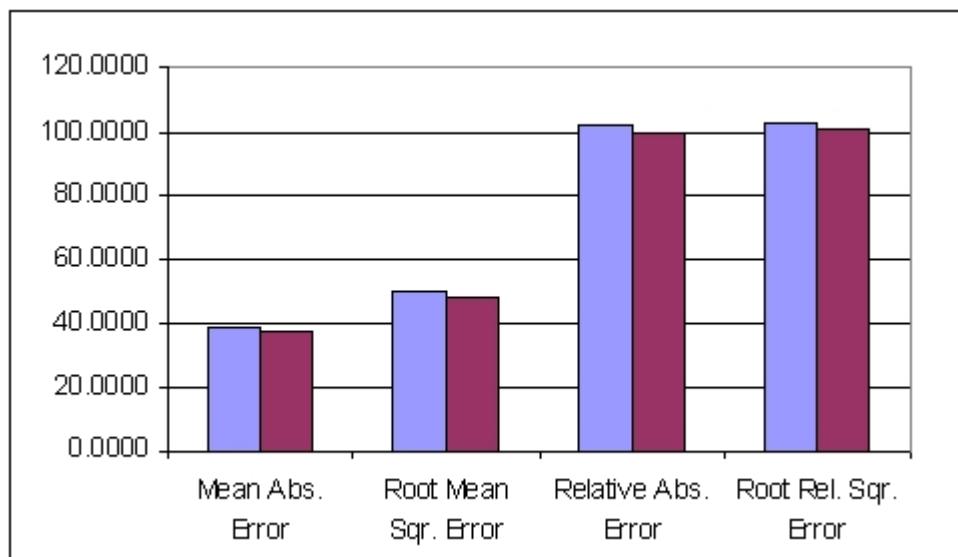
c)



d)



e)

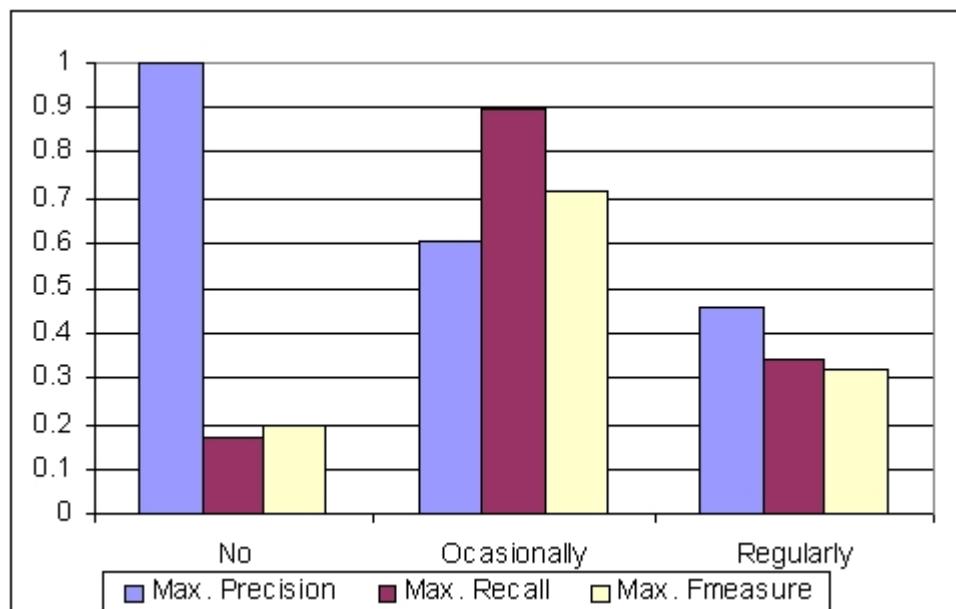


f)

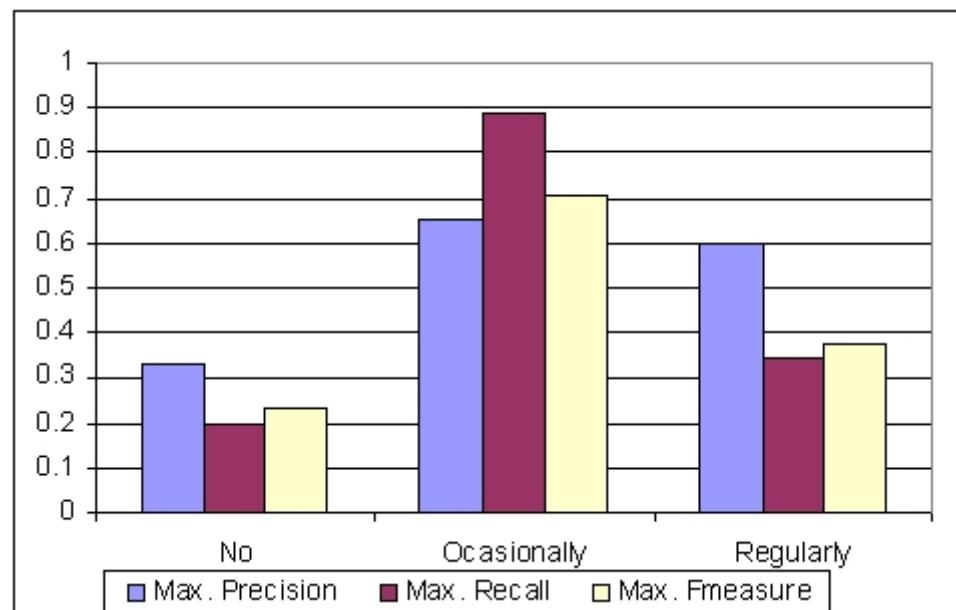
*Figura 6. Media y máximo del error absoluto, error cuadrático medio, error relativo y error cuadrático relativo de todos los algoritmos para la predicción del nivel de colesterol, dados los factores sociales en a) grupo Normal, b) grupo Patológico y c) grupo de Riesgo, y dada la actividad física en d) grupo Normal, e) grupo Patológico y f) grupo de Riesgo.*

Veamos otro ejemplo representativo. La Figura 6 presenta los resultados de la predicción del nivel de colesterol dado los factores sociales para, a), b) y c), y dada la actividad física, d), e) y f) para cada grupo de nivel respectivamente. En este caso, los resultados de predicción son muy similares para los dos conjuntos de atributos de entrenamiento en todos los grupos de nivel, por lo que podemos concluir que las fuerzas de las relaciones también lo son. Sin embargo, la predicción varía de un grupo de nivel a otro. En el grupo Normal, el error de predicción medio es de alrededor de 24, siendo entorno a 50 y 40 en los grupos Patológico y Normal, respectivamente, concluyendo que es más fácil predecir el nivel de colesterol, ya sea a partir de factores sociales o de actividad física, para los sujetos en el grupo Normal. Este hecho denota una fuerte relación entre los factores de entrenamiento y el nivel de colesterol en este grupo de nivel.

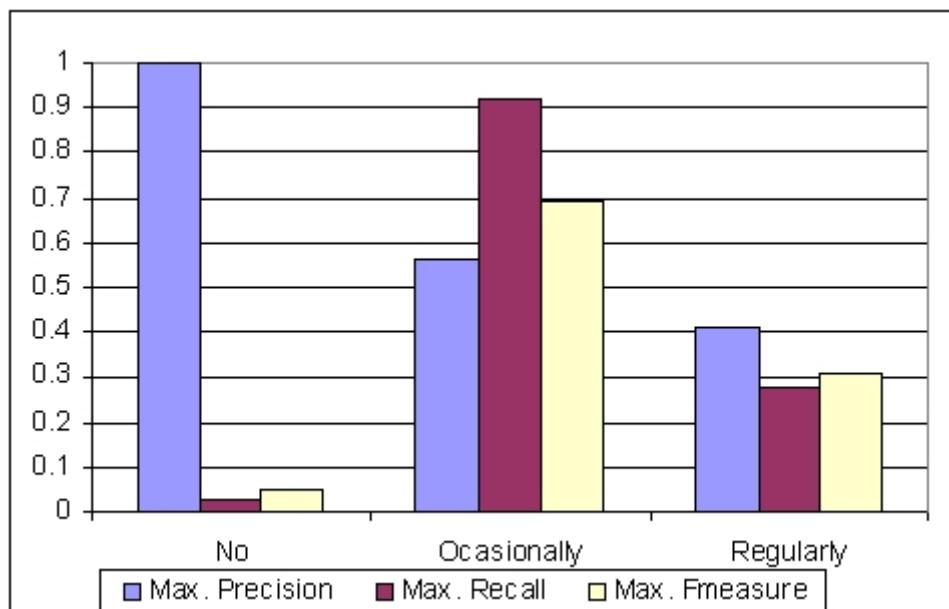
Finalmente, la Figura 7 muestra los resultados de predicción de los valores del atributo “Alcohol”, dados los factores sociales y la actividad física por separado para cada uno de los grupos de nivel, análogamente a las figuras anteriores.



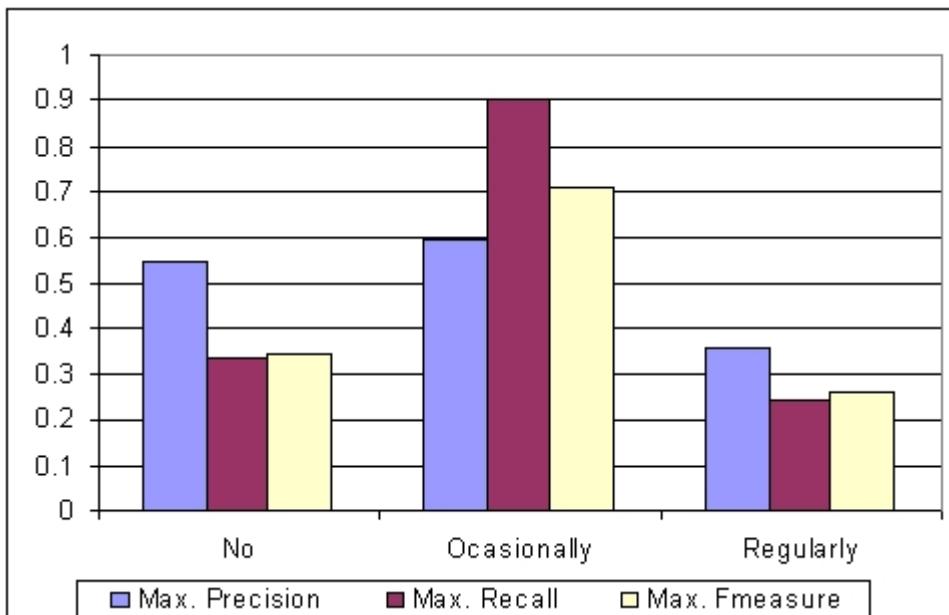
a)



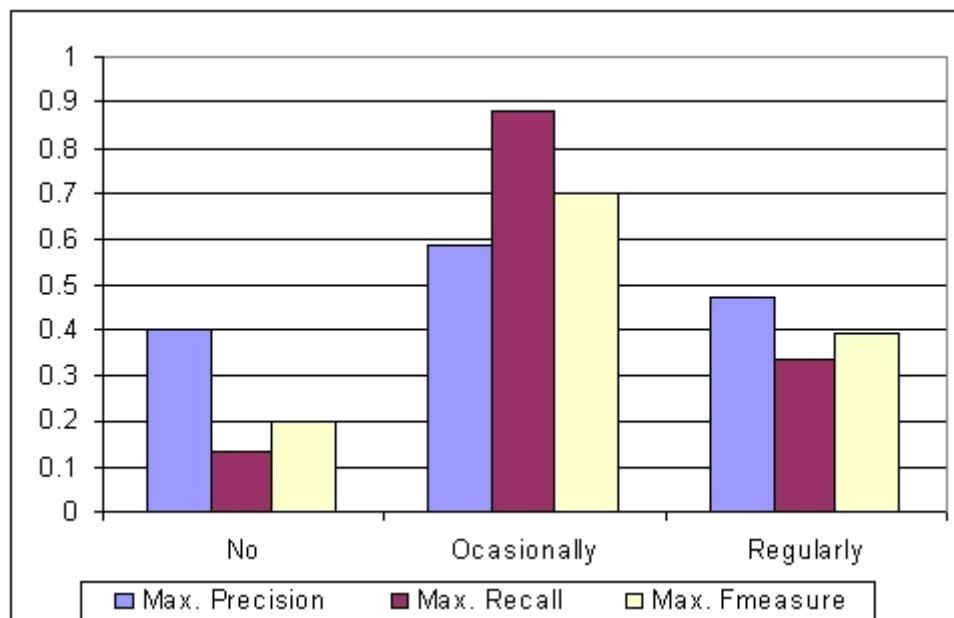
b)



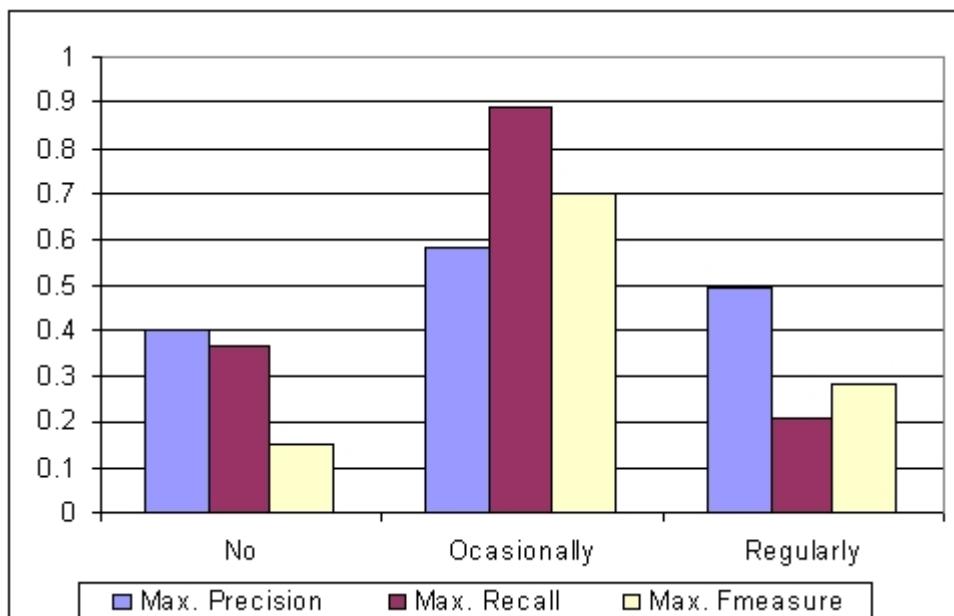
c)



d)



e)



f)

Figura 7. Precisión, cobertura y medida-F máximas de entre todos los algoritmos para la predicción de "Alcohol", dados los factores sociales en a) grupo Normal, b) grupo Patológico y c) grupo de Riesgo, y dada la actividad física en d) grupo Normal, e) grupo Patológico y f) grupo de Riesgo.

Los resultados de la Figura 7 muestran que existe una clara relación en todos los grupos de nivel entre los atributos de entrenamiento y los sujetos que beben alcohol ocasionalmente. Los sujetos que beben de forma regular son más difíciles de ser detectados a partir de los factores de entrenamiento, presentando una leve asociación ligeramente más significativa en el grupo

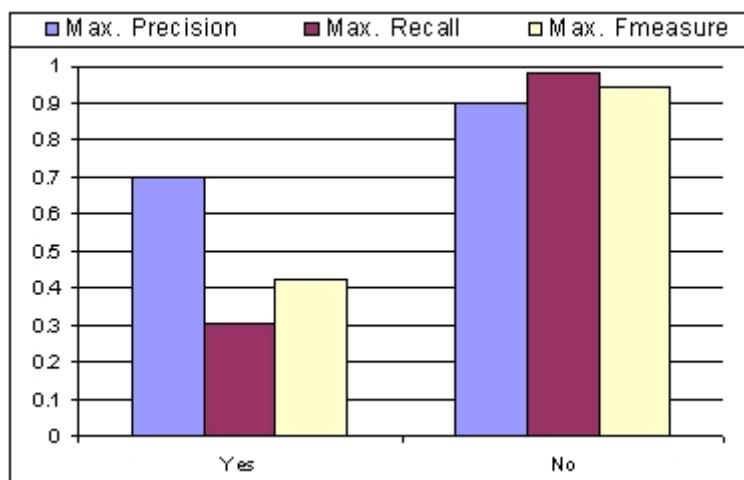
Patológico. Se puede decir lo mismo, con respecto a la actividad física, para los sujetos que nunca beben alcohol. Sin embargo, la precisión de la predicción se incrementa notablemente dados los factores sociales en los grupos Normal y de Riesgo. Los sujetos que nunca beben son identificados de manera precisa a partir de los atributos sociales, lo que denota una relación significativa entre los factores involucrados.

Los atributos de entrenamiento se presentan como entrada conjunta a los algoritmos. Desde el punto de vista médico es también interesante separar estos atributos y presentar combinaciones de los mismos. Así, se intentó predecir el valor de actividad física en el trabajo a partir de todas las posibles combinaciones de atributos sociales. Los resultados mostraron que, para los grupos Normal y de Riesgo, el atributo de entrenamiento de nivel educativo por sí solo obtenía resultados mucho mejores que cualquiera del resto de posibles combinaciones de atributos sociales. En el grupo Patológico los resultados son similares, aunque la diferencia no es tan notoria como en los otros grupos, siendo Edad+Nivel Educativo la mejor combinación de atributos para predecir la actividad física en el trabajo.

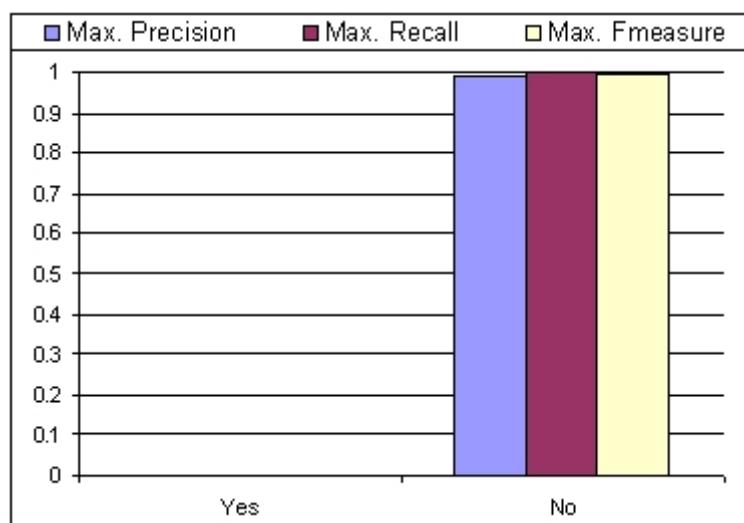
## 5.2 Predicción de disfunciones futuras

El objetivo principal de los experimentos descritos a continuación es comprobar la eficacia de la predicción de disfunciones futuras. En este caso, además del conjunto *Entry*, el conjunto *Control* de la colección también se emplea. Primeramente, se seleccionaron los pacientes con registros posteriores a 10 años desde su entrada en el estudio del subconjunto *Control*. Despues, usando sus atributos correspondientes del conjunto *Entry* se intentó predecir si padecerían alguna enfermedad en 10 años. Las enfermedades consideradas son hipertensión sistólica, diastólica o sistólica-diastólica, hipercolesterolemia y hipertrigliceridemia. Los valores correspondientes a estos atributos son verdadero o falso, es decir, la padecen o no. La misma tarea de predicción también fue realizada para disfunciones en 20 años. Los resultados muestran que el Perceptrón Multicapa fue el mejor algoritmo, alcanzando valores cercanos al 85% de precisión y 65% de cobertura en la detección de todas las enfermedades. Puesto que el riesgo de hipertensión en el grupo de Riesgo es nulo y algunos de estos pacientes padecían hipertensión desde el principio del estudio, es más interesante desde el punto de vista médico la predicción de enfermedades en el grupo Normal. Así pues, el mismo proceso se ha realizado sólo en el grupo mencionado para la predicción a 10 y a 20 años. Los resultados para las diferentes disfunciones se presentan en la Figura 8, a) a e), respectivamente para 10 años y f) a j), respectivamente para 20 años. Para cada disfunción, se muestran los valores máximos de entre todos los algoritmos. En este caso no existe ningún algoritmo que destaque con respecto a los demás. Dependiendo de la disfunción algunos algoritmos funcionan mejor que otros, por lo que es interesante usar todos los métodos y tomar decisiones a partir de sus resultados conjuntos. Es interesante comentar que la predicción es mucho más eficaz cuando se consideran todos los grupos de nivel conjuntamente, confirmando el interés previo en el grupo Normal.

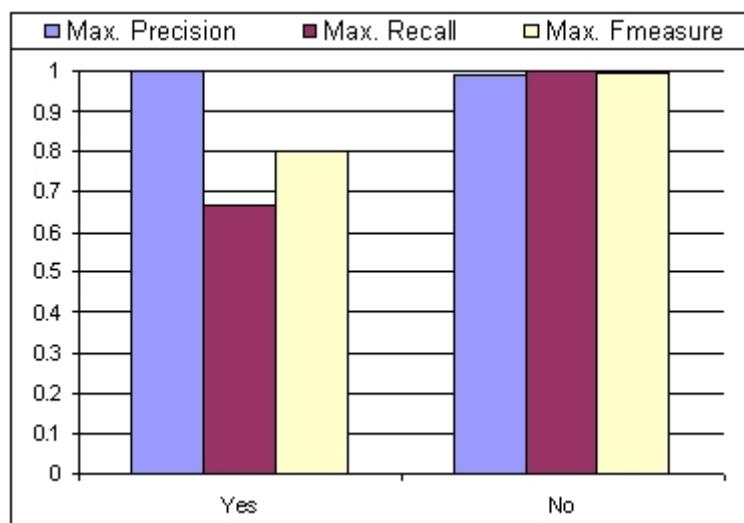
Los valores de la Figura 8 muestran que es más precisa la predicción a 20 años que a 10 años. De hecho, la detección de la presencia (y no la ausencia) de disfunciones se infiere de manera más precisa a 20 años. La ausencia de disfunciones es igualmente bien inferida para cualquier intervalo de tiempo. Entre todas las disfunciones la mejor detectada en la hipertensión diastólica, obteniendo valores entorno al 100% de precisión tanto para la ausencia como para la presencia de la disfunción. La más difícil de detectar es la hipertensión diastólica, imposible de detectar en 10 años.



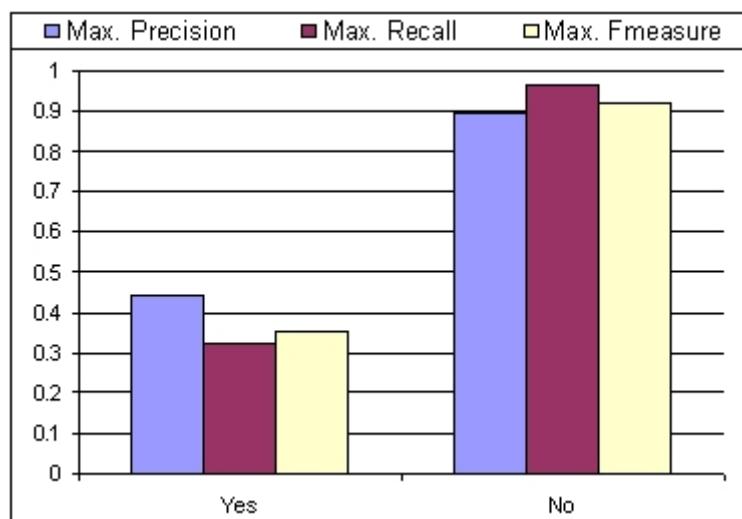
a)



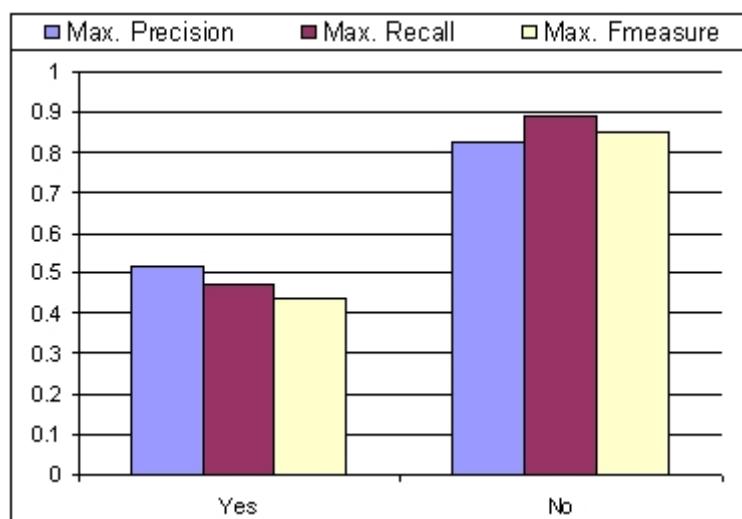
b)



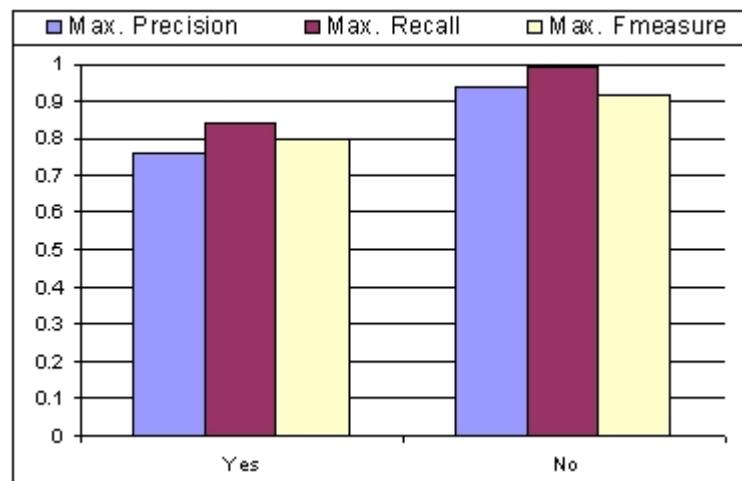
c)



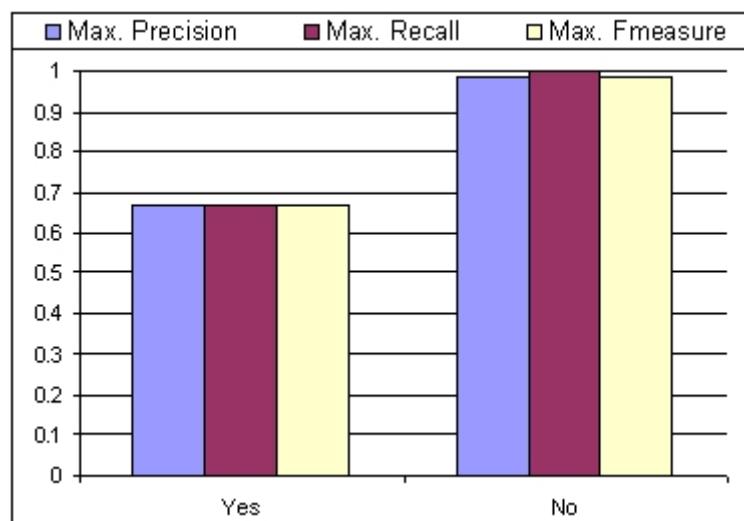
d)



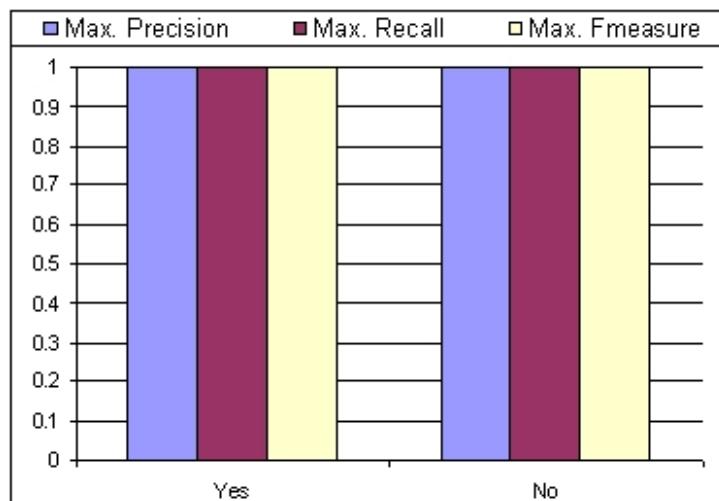
e)



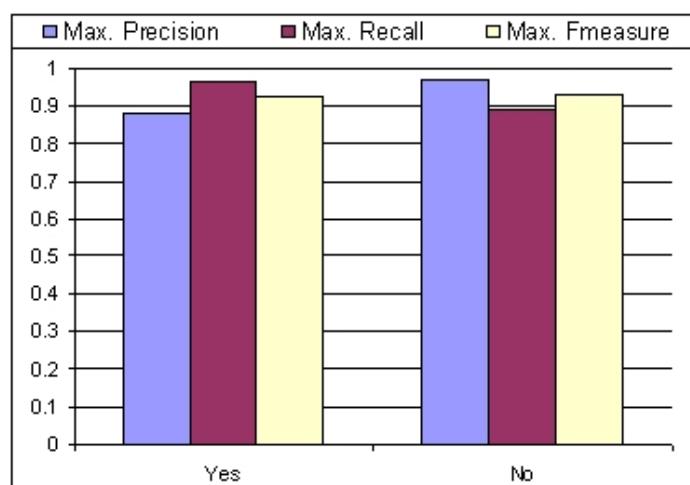
f)



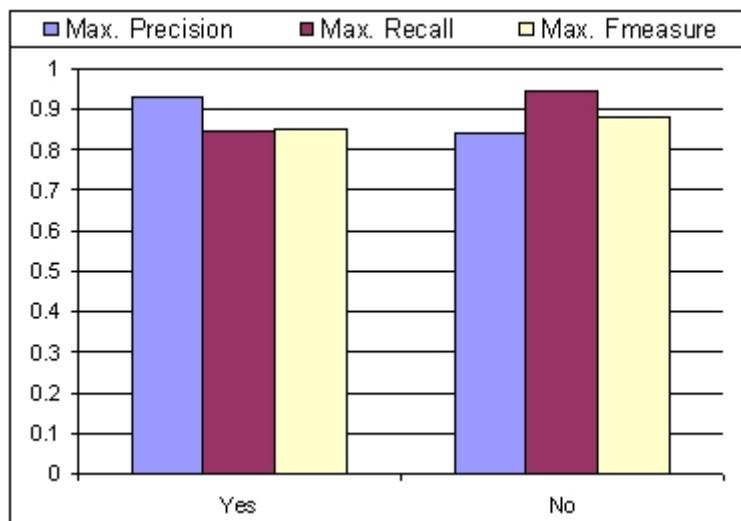
g)



h)



i)



j)

*Figura 8. Máxima precisión, cobertura y medida-F para la predicción de a) hipertensión sistólica-diastólica, b) hipertensión sistólica, c) hipertensión diastólica, d) hipercolesterolemia y e) hipertrigliceridemia en 10 años, y f) hipertensión sistólica-diastólica, g) hipertensión sistólica, h) hipertensión diastólica, i) hipercolesterolemia y j) hipertrigliceridemia en 20 años.*

También se intentaron predecir otras disfunciones como la angina de pecho, el infarto de miocardio, infarto cerebro-vascular, etc., pero dado el escaso número de ejemplos correspondiente no se puede extraer ninguna conclusión de los resultados.

### 5.3 Predicción de la causa de muerte

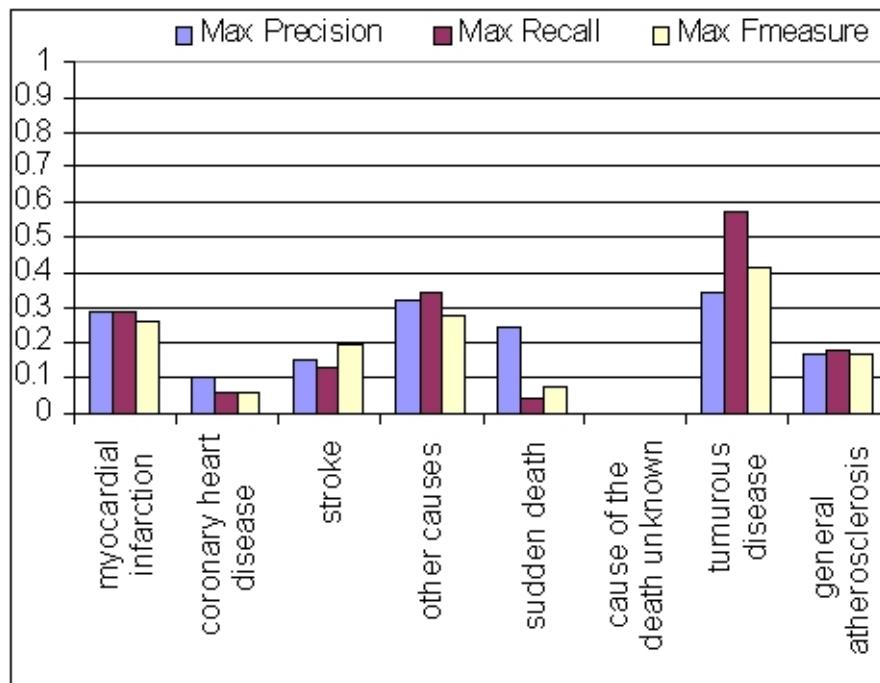
Este experimento es análogo al anterior, pero ahora se trata de predecir la causa de muerte, por lo que se emplea el conjunto *Death* de la colección. Los algoritmos fueron entrenados con los datos del conjunto *Entry* de los pacientes presentes en el conjunto *Death*. Los experimentos se realizaron para cada grupo de nivel por separado y para los tres juntos. Los resultados se presentan en la Figura 9. En el grupo Normal, Figura 9 b), la causa de muerte más detectada fue la enfermedad tumoral. En el grupo de Riesgo, Figura 9 d), otras causas desconocidas, además del infarto de miocardio y enfermedad coronaria fueron las causas mejor inferidas, en absoluto detectadas en el grupo Normal. En el grupo Patológico, Figura 9 c), las causas mejor inferidas fueron la enfermedad tumoral y el infarto de miocardio. Aunque el paro cardíaco y la arterosclerosis general fueron levemente detectadas en el grupo Patológico, en los anteriores grupos no pudieron ser inferidas. De manera general con todos los grupos unidos, Figura 9 a), los resultados de predicción son muy pobres, concluyendo que los datos del conjunto *Entry* no poseen suficiente información para poder predecir la causa de muerte, o quizás se necesiten más registros de muertes de pacientes.

## 6. Conclusiones

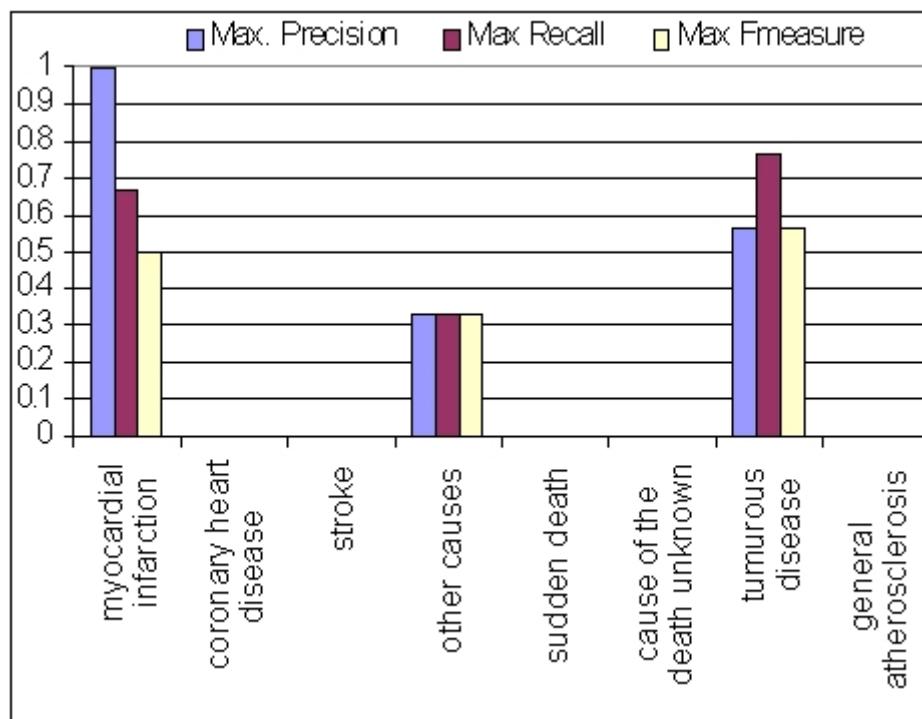
Se han aplicado diferentes algoritmos de aprendizaje automático al descubrimiento de conocimiento en datos biomédicos de dos maneras distintas: primero, los métodos se han

Métodos de Aprendizaje Automático para el Descubrimiento de Conocimiento en Datos Médicos...

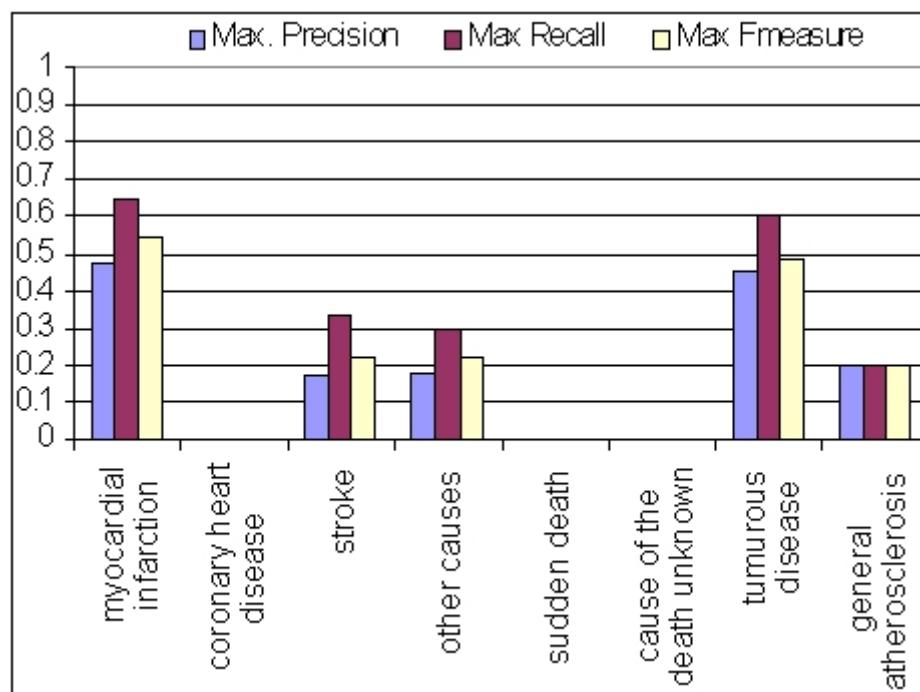
usado para predecir el valor de un atributo de la colección dado un subconjunto de otros atributos como entrenamiento, proponiendo la máxima eficacia entre todos los algoritmos como medida de la fuerza de la asociación entre los atributos de entrenamiento y el atributo objetivo. Esta medida ha resultado útil también para la comparación de las asociaciones en diferentes grupos de pacientes.



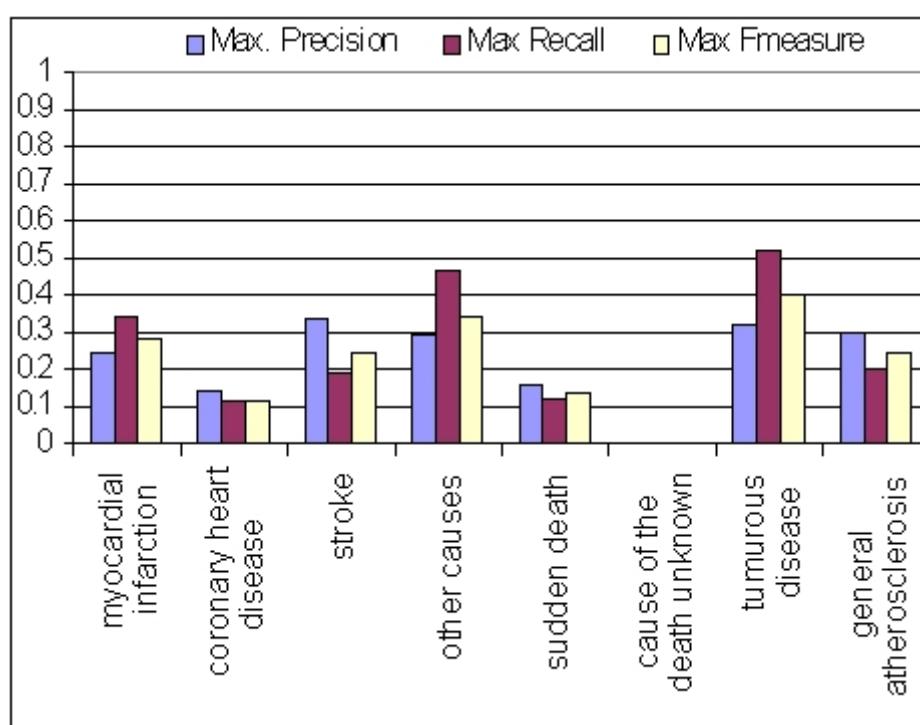
a)



b)



c)



d)

Figura 9. Máximos valores de precisión, cobertura y medida-F en la predicción de causa de muerte para a) todos los grupos de nivel juntos, b) grupo Normal, c) grupo Patológico y d) grupo de Riesgo.

En segundo lugar, las técnicas de aprendizaje se han aplicado a la predicción de disfunciones futuras. Los resultados muestran que algunos métodos predicen ciertos desórdenes mejor que otros, por lo que es interesante usar todos los algoritmos y considerar sus resultados en base a la tendencia conocida de cada método. Todos los algoritmos empleados predicen mejor a 20 años que a 10, alcanzando excelentes resultados para algunas de las disfunciones, lo que les hace adecuados para la ayuda a la toma de decisiones. Los algoritmos también han sido evaluados en la predicción de causas de muerte, obteniendo resultados poco significativos debido quizás a la escasa información de este tipo presente en la colección.

En un futuro cercano se pretende ajustar y optimizar los parámetros de los algoritmos y evaluar más métodos. Se pretende también integrar todos los algoritmos con los grados de significación y utilidad hallados en este trabajo para construir un sistema experto. Se investigará también la derivación de reglas a partir de los resultados de los algoritmos que sean interpretables por los médicos y especialistas.

## Agradecimientos

Esta investigación ha sido conjuntamente financiada por el Plan de Investigación de ICS AS CR AV0Z10300504 y por el Plan de Estancias Breves en el Extranjero “María Bueno” del Consejo Superior de Investigaciones Científicas junto con el Instituto de Automática Industrial del CSIC.

## References

- [1] Mitchell, T.: Machine Learning. McGraw Hill, 1997.
- [2] Lavrać, N.: Selected Techniques for Data Mining in Medicine. Artificial Intelligence in Medicine, vol. 16 (1), pp. 3-23, 1999.
- [3] Aseervatham, S. and Osmani A.: Mining Short Sequential Patterns for Hepatitis Type Detection. ECML/PKDD Discovery Challenge, 2005.
- [4] Aubrecht, P., Kejkula, M., Kremen, P., Novakova, L., Rauch, J., Simunek, M., Stepankova, O.: Mining in Hepatitis Data by LISp-Miner and SumatraTT. ECML/PKDD Discovery Challenge, 2005.
- [5] Pizzi, L.C., Ribeiro, M.X., Vieira, M.T.P.: Analysis of Hepatitis Dataset using Multirelational Association Rules. ECML/PKDD Discovery Challenge, 2005.
- [6] Durand, N., Soulet, A.: Emerging Overlapping Clusters for Characterizing the Stage of Liver Fibrosis. ECML/PKDD Discovery Challenge, 2005.
- [7] Durand, N., Cleuziou, G., Soulet, A.: Discovery of Overlapping Clusters to Detect Atherosclerosis Risk Factors. ECML/PKDD Discovery Challenge, 2004.
- [8] Cios, K. J.: Medical data mining and Knowledge Discovery. Physica – Verlag, 2001.
- [9] Chen, H., Fuller, S. S., Friedman, C. and Hersh, W.: Medical Informatics: Knowledge Management and Data Mining in Biomedicine. Integrated Series in Information Systems (2), Springer Science and Business Media Inc., 2005.
- [10] Boudik F., Reissigova J., Hrach K., Tomeckova M., Bultas J., Anger Z., Aschermann M., Zvarova J.: Primary Prevention of Coronary Artery Disease Among Middle Aged Men in Prague: Twenty-year Follow-up Results. Atherosclerosis. 2006 Jan;184(1):86-93.

- [11] Tomeckova, M.: The Challenge on Atherosclerosis Data Viewed by the Experts. ECML/PKDD Discovery Challenge, 2004.
- [12] Rish, I.: An Empirical Study of the Naive Bayes Classifier. IJCAI-01 Workshop on Empirical Methods in AI, 2001.
- [13] Haykin, S.: Neural Networks: A comprehensive Foundation (2nd edition). Pearson Education, 1998.
- [14] Scholkopf, B., Smola, A. J., Mtiller, K.-R., Burges, C. J. C., and Vapnik, V.: Support Vector Methods in Learning and Feature Extraction. In Down, T., Frean, M., and Gallagher, M., editors. Proceedings of the Ninth Australian Congress on Neural Networks, Brisbane, Australia. University of Queensland, 1998.
- [15] Teknomo, K.: K-Nearest Neighbors Tutorial. <http://people.revoledu.com/kardi/tutorialKNN>, 2004.
- [16] Quinlan, J. R.: C4.5: Programs for Machine Learning. Morgan Kauffman, 1993.
- [17] Compton, P., Edwards, G., Kang, B., Malor, R., Menzies, T., Preston, P., Srinivasan, A. and Sammut, S.: Ripple Down Rules: Possibilities and Limitations. Boose, J.H. & Gaines, B.R., Ed. Proceedings of the Sixth AAAI Knowledge Acquisition for Knowledge-Based Systems Workshop. pp.6-1-6-20. Calgary, Canada, University of Calgary, 1991.
- [18] Van Rijsbergen, C. J.: Information Retrieval. Butterworths, London, 1979.
- [19] Witten, I. H. and Frank, E.: Data Mining: Practical Machine Learning Tools and Techniques. 2nd Edition, Morgan Kaufmann, San Francisco, 2005.

# Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze

José Ignacio Serrano<sup>1</sup>, Marie Tomečková<sup>2</sup>, Jana Zvárová<sup>2</sup>

1. Instituto de Automática Industrial, CSIC, Madrid, Spain,

2. Department of Medical Informatics, Institute of Computer Science AS CR, Prague, Czech Republic

Techniky strojového učení jsou metody, které umožní vytvořit z trénovací množiny případů model pro kategorie dat tak, že mohou být nové (neznámé) případy zařazeny do jedné nebo více kategorií schématem odpovídajícím modelu. Pro tento typ analýzy jsou velmi vhodná data ze studií sledujících určitou skupinu osob s opakovaným sběrem dat stejného typu. K vyhledávání znalostí z medicínských dat bylo užito různých algoritmů strojového učení. Bylo testováno několik algoritmů tak, aby bylo možno pokrýt většinu způsobů učení s učitelem. Byly provedeny dva typy pokusů. Jeden hledal vztahy mezi atributy, druhý testoval predikci budoucích příhod. Pro pokusy v tomto sdělení byla užita data z dvacet let trvající longitudinální primárně preventivní studie rizikových faktorů (RF) aterosklerózy u mužů středního věku. Studie se nazývá STULONG (LONGitudinal STUDy). Výsledky ukazují, že některé metody předpovídají některé poruchy lépe než jiné a že je tedy vhodné použít všechny algoritmy najednou a posuzovat spolehlivost výsledku na základě známého trendu každé metody. Algoritmy strojového učení byly také použity k předpovědi příčiny úmrtí. V tomto případě byly výsledky nevalné, pravděpodobně pro malé množství informace ve vstupních položkách v datového souboru.

**Klíčová slova:** dobývání znalostí, strojové učení s učitelem, vytěžování z biomedicínských dat, rizikové faktory aterosklerózy

## 1. Úvod

Techniky strojového učení [1] jsou metody, které umožní vytvořit z trénovací množiny případů model pro kategorie dat, takže mohou být nové (neznámé) případy zařazeny do jedné nebo více kategorií schématem odpovídajícím modelu.

Techniky strojového učení byly úspěšně použity k řešení předpovědních úloh u řady různých problémů a dat. Základním úkolem je zjistit, jak použít algoritmus strojového učení na tato data s cílem odhalit vztahy mezi atributy a vytvořit predikce, které by mohly být užitečné pro podporu rozhodování. Medicínská data jsou specifický druh dat, protože při jejich sběru je zaznamenáváno mnoho různých druhů atributů. Nicméně, medicínská data mají několik známých problémů: chybějící, nesprávné nebo málo četné informace a časově omezená data. Pro tento typ dat jsou velmi vhodné metody strojového učení [2]. Řada prací KDD (knowledge discovery from databases – získávání znalostí z databází) se snaží pracovat s rozsáhlým množstvím lékařských informací. V práci [3] se autoři pokouší určit typ zánětu jater tím, že vyberou krátký sled charakteristik z časově ohraničeného záznamu. V práci [4] byla podobná úloha řešena pomocí metody čtyřpolních tabulek (tj. statistických tabulek

s dvěma řádky a dvěma sloupcí), aby bylo možno stanovit zánět jater typu B a C z rozdílů onemocnění. Autoři v [5] se pokusili řešit úlohu jednoduchými boolen údaji (ano-ne), které mohou předpovídat stadium jaterní cirkózy. Podobné aplikace užili autoři v [6], ale v tomto případě byly vybrané příznaky spojeny a tyto bloky příznaků byly přiřazeny ke stadiu jaterní cirkózy, což záviselo na zahrnutých případech. Tato technika byla také použita při určování rizika aterosklerózy [7]. Pro analýzu jsou velmi vhodné studie se sledováním osob a opakováním shromažďováním dat stejného typu. Další příklady získávání znalostí z biomedicínských dat ukazují práce [8] a [9]. Pro pokusy v tomto sdělení byla použita data z dvacetileté longitudinální primárně preventivní studie rizikových faktorů (RF) aterosklerózy u mužů středního věku. Studie se nazývá STULONG (LONGitudinální STUDie) [10], [11]. Hlavním cílem těchto pokusů je vyhodnotit strojové učení jako způsob vyhledávání asociací a zhodnotit výstup klasifikace pro měření charakteristických rysů nalezených asociací. Algoritmy strojového učení jsou také používány pro testování pro předpověď poruch ve vzdálené budoucnosti.

V následující části jsou uvedeny podrobnosti o souboru dat studie STULONG. V části 3 jsou popsány testované algoritmy strojového učení. Část 4 předkládá měření pro hodnocení a Část 5 popisuje validaci pokusů. Na závěr jsou v části 6 uvedeny závěrečné poznámky a návrh další práce.

## 2. Popis studie a souboru dat

Data studie STULONG (<http://euromise.vse.cz/challenge2004/index.html>) [10], [11] byla získávána v letech 1975-1999 na II. interní klinice 1. lékařské fakulty Univerzity Karlovy v Praze a Všeobecné fakultní nemocnice v Praze. Data byla převáděna do elektronické formy Evropským centrem pro lékařskou informatiku, statistiku a epidemiologii Univerzity Karlovy v Praze a Akademie věd České republiky v Praze v letech 1994-1999 a poté analyzována statisticky.

Hlavní cíle studie byly:

1. Zjistit prevalenci rizikových faktorů (RF) aterosklerózy v populaci, která je všeobecně považována za nejohroženější možnými komplikacemi aterosklerózy, tj. u mužů středního věku.
2. Sledovat vývoj těchto RF a jejich dopad na zdraví vyšetřených mužů, zejména na vznik aterosklerotických srdečně-cévních onemocnění.
3. Zhodnotit dopad komplexní intervence RF na vývoj těchto RF a na srdečně-cévní nemocnost a úmrtnost sledovaných mužů.

V roce 1975 byli z volebních seznamů Prahy 2 vybráni muži narození v letech 1926-1937 a žijící v obvodu Prahy 2. Na první vyšetření se dostavilo z 2370 pozvaných 1419 mužů. Vstupní vyšetření byla provedena v letech 1976-1979. V dopisu, kterým byli muži zváni k vyšetření, byly krátce vysvětleny cíle studie, průběh vstupního vyšetření a účel dalšího sledování. Muži byli požádáni o spolupráci. V té době nebyl vyžadován k účasti na studii podpis informovaného souhlasu. Pokud muž reagoval na pozvání a přišel na první vyšetření, považovali jsme to za dostatečný souhlas s vyšetřením, dalším sledováním a analýzou výsledků. Pokud muž na první pozvání nereagoval, zaslali jsme mu další, maximálně dvě, pozvání.

Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze

Rizikové faktory byly stanoveny podle tehdejších definic následovně:

- hypertenze – krevní tlak  $\geq 160/95$  mm Hg nebo muži, užívající léky ke snížení krevního tlaku,
- hypercholesterolemie – celkový cholesterol  $\geq 260\text{mg\%}$  (6,7 mmol/l),
- hypertriglyceridemie – triglyceridy  $\geq 200\text{mg\%}$  (2,2 mmol/l),
- kouření:  $\geq 15$  cigaret/den v současné době nebo stejně množství cigaret kouřených denně v době kratší než 1 rok před vstupem do studie (kuřáci dýmek nebo doutníků byli zařazeni mezi nekuřáky),
- nadváha: Brocka index  $> 115\text{\%}$  (Brocka index: výška v cm - 100 = 100 %),
- pozitivní rodinná anamnéza: úmrtí otce nebo matky na ischemickou chorobu srdeční nebo cévní mozkovou příhodu před jejich 65. rokem věku.

Podle přítomnosti RF, celkového zdravotního stavu a nálezu na záznamu EKG byli muži rozděleni do následujících skupin:

- NG** = Normální skupina – skupina mužů bez RF definovaných výše, bez manifestního aterosklerotického onemocnění nebo jiného závažného onemocnění, které by bránilo jejich desetiletému sledování a beze změn na EKG záznamu.
- RG** = Riziková skupina – skupina mužů s alespoň jedním RF podle výše uvedené definice, bez manifestního aterosklerotického onemocnění nebo jiného závažného onemocnění, které by bránilo jejich desetiletému sledování a beze změn na EKG záznamu.
- PG** = patologická skupina – skupina mužů s manifestním aterosklerotickým onemocněním nebo jiným závažným onemocněním, které by bránilo jejich desetiletému sledování (např. maligní onemocnění, pokročilé jaterní nebo ledvinné selhání, závažné neurologické nebo psychické poruchy). V patologické skupině byli též muži s cukrovkou, léčení perorálními antidiabetiky nebo inzulínem a muži s patologickým nálezem na EKG záznamu podle Minnesotského kódu.

Pro dlouhodobé sledování byli pacienti rozděleni do následujících skupin:

- Riziková skupina **RG** byla randomizovaně rozdělena na dvě podskupiny označené jako **RGI** (intervenovaná riziková skupina) a **RGC** (kontrolní riziková skupina). Pacienti **RGI** skupiny byli zváni na kontrolu nejméně dvakrát ročně, při farmakologické intervenci byli zváni podle potřeby. Pacienti z **RGC** skupiny obdrželi krátkou písemnou zprávu včetně laboratorních výsledků a popisu EKG záznamu s doporučením, aby tyto výsledky předali svému ošetřujícímu lékaři. Případná intervence RF byla ponechána na rozhodnutí tohoto lékaře. Muži ze skupiny RGC byli zváni na kontrolu jedenkrát ročně. Při prvním vyšetření nebyly mezi skupinami RGI a RGC signifikantní rozdíly ve věku, socio-ekonomických ukazatelích, ani ve výskytu RF.
- 10 % mužů z **NG** skupiny bylo kontrolováno jedenkrát ročně podobně muži rizikových skupin (byli označeni jako NGS – sledovaní). V této skupině – podobně jako ve skupině RGI – byla zahájena intervence RF co nejdříve po jejich zjištění (hypertenze, hyperlipidemie). Ostatní muži z NG byli pozváni na kontrolu za 10-12 let po prvním vyšetření.
- Muži z **PG** skupiny byli ze sledování vyloučeni.

## Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze

Intervence RF měla ve studii klíčové postavení a byla vždy zahájena nefarmakologickými opatřeními. Snažili jsme se o úpravu RF a dosažení jejich optimálních hodnot.

- *Nefarmakologická intervence*: pohovor o životním tylu, tj. stravování, tělesné aktivitě, vhodnosti, resp. nutnosti zanechat kouření a snížit váhu. Tyto pohovory byly při každé kontrole opakovány a kromě obecných doporučení byly zaměřeny na RF u daného pacienta.
- *Farmakologická intervence*: léčba arteriální hypertenze a hyperlipoproteinemie – v počátku studie byly velmi omezené a v širším měřítku byly užity až v posledních letech studie. Farmakologická léčba byla doporučena podle celkového rizika a případných dalších onemocnění pacienta.

K analýze byly použity čtyři soubory dat:

1. Soubor *ENTRY (Vstup)* obsahuje 244 atributů ze vstupního vyšetření každého muže; tyto atributy jsou výsledky různých veličin, často kódované, nebo jsou výsledkem transformací původních veličin (identifikace muže, rodinná a osobní anamnéza, sociální faktory – vzdělání, tělesná aktivita, kouření, stravovací návyky, spotřeba alkoholu, poté antropometrická měření – výška, váha, kožní řasy, fyzikální vyšetření včetně změření krevního tlaku, zjištění tepové frekvence, laboratorní hodnoty a kódovaný EKG záznam).
2. Soubor *CONTROL (Kontroly)* obsahuje 66 atributů zaznamenaných při kontrolním vyšetření. Tyto atributy odpovídají identifikaci muže, záznamům o změnách ve způsobu života, v osobní anamnéze, fyzikální vyšetření, výsledky biochemického vyšetření a údaje o hypertenzi, hypercholesterolemii, hypertriglyceridemii a prodělaných onemocněních, zejména srdečních a nádorových. Tento soubor obsahuje 10 572 záznamů z dlouhodobého sledování.
3. Dodatečné informace o zdravotním stavu 403 mužů, kteří ukončili studii předčasně, byly získány dotazníkem zaslaným mužům poštou. Bylo získáno 62 atributů a jsou v souboru *LETTER (Dopis)*.
4. V souboru *DEATH (Úmrtí)* je 5 atributů o každém z 389 pacientů, kteří během studie zemřeli. Tyto atributy jsou identifikace pacienta a datum a příčina úmrtí.

## 3. Popis použitých metod

Všechny použité algoritmy patří mezi učící postupy s učitelem. To znamená, že je třeba mít učící množinu k vytvoření modelu trénovacích případů a poté použít tento model k předpovědi kategorie neznámých případů. Bylo testováno několik algoritmů ve snaze podchytit řadu způsobů učení s učitelem. V následující části je velmi krátce vysvětlena každá z těchto metod.

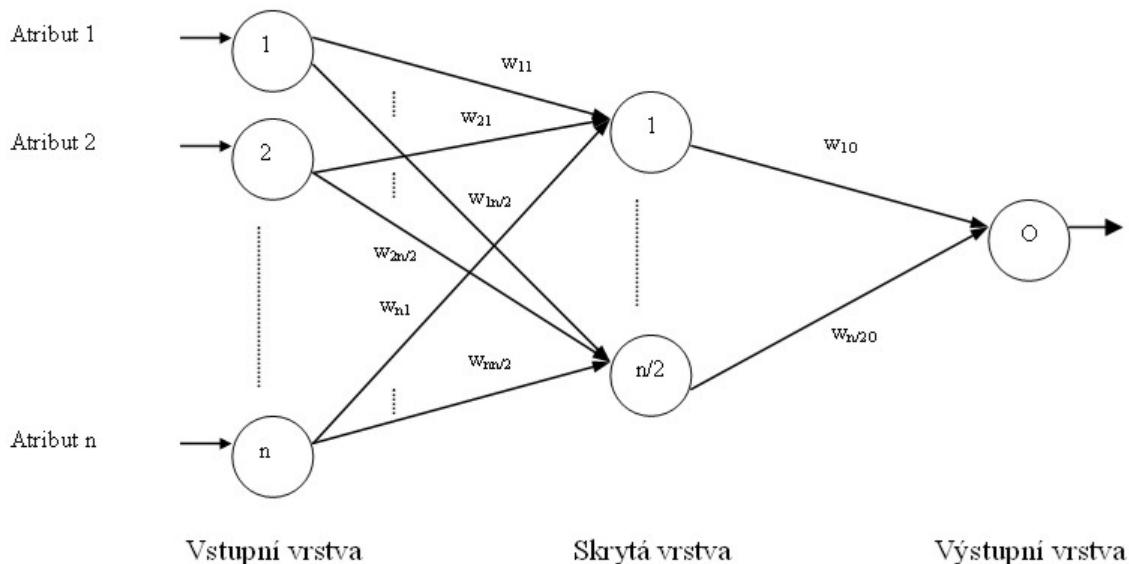
### 3.1 Naivní Bayes

Naivní Bayes [11] počítá pro každý pár „atribut a jeho hodnota“ (např. vzdělání, vysokoškolské) pravděpodobnost příslušnosti ke každé kategorii a to tak, že dělí počet případů v dané kategorii, kde se pár vyskytuje, celkovým počtem případů v celém souboru, kde se pár vyskytuje. Každý pár bude mít určitou pravděpodobnost pro každou uvažovanou kategorii. Metoda je založena na předpokladu, že každý pár „atribut-hodnota“ je nezávislý na

jakémkoliv jiném páru. Takže, když je klasifikován neznámý případ, pravděpodobnost příslušnosti ke každé kategorii je násobkem pravděpodobností každého páru, který tvoří případ, pro odpovídající kategorii. Předpověděná kategorie je ta s největší pravděpodobností.

### 3.2 Vícevrstevný perceptron

Klasifikační model neuronové sítě vícevrstevného perceptronu [13] je tvořen sítí složenou z vrstev vzájemně propojených neuronů. Mezi neurony jedné vrstvy nejsou žádné vazby, ale neuron z jedné vrstvy je propojen se všemi neurony vrstvy sousední. Architektura modelu užitého v našem souboru je ukázána na Grafu 1.



Graf 1. Architektura užité neuronové sítě vícevrstevného perceptronu.

S každým spojením je asociována určitá váha. Vstup (input) do každého neuronu je vážený součet spočtený z asociovaných vah všech přicházejících hodnot. Výstup z každého neuronu je výsledek použité funkce. V tomto případě je u všech neuronů použita typická sigmoidní funkce. Graf 2 ukazuje vyjádření a zobrazení této funkce.



Graf 2. Vyjádření a zobrazení esovité funkce.

Každá hodnota atributu ze vzorku z datového souboru je vložena do odpovídajícího neuronu ve vstupní vrstvě a hodnoty procházejí neuronovou sítí k výstupní vrstvě, kde výstupní hodnota neuronu znamená předpovídanou kategorii.

Trénovací fáze spočívá ve vložení každého označeného případu z trénovacího datového souboru s daným souhrnem původních vážených hodnot do modelu a porovnání výstupní hodnoty s předpokládanou kategorií. Podle chyby v předpověděné kategorii jsou při zpětném šíření algoritmu změněny váhy z výstupní vrstvy do vstupní vrstvy tak, aby předpovídána hodnota byla předpokládané hodnotě podobnější. Tento proces je prováděn s určitým počtem iterací. V našem případě je tento počet 500. Rozsah, o který se mění váhy při zpětném šíření, tzv. učící poměr, je 0,3. Pokud algoritmus zpětného šíření nedosahuje dobrého přiblížení k předpokládanému výstupu po jedné iteraci, potom se model obnoví a směřuje ke snížení učícího poměru.

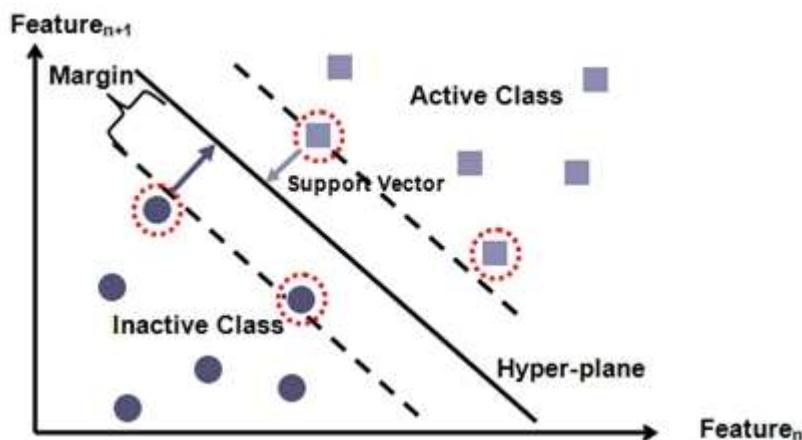
### 3.3 Support Vector Machines (SVM)

SVN (Support Vector Machines) [14] se pokouší rozdělit případy založené na jejich kategoriích do  $n$ -rozměrného prostoru;  $n$  je počet atributů nebo charakteristik, nadrovina má vyjádření  $w + b$ , takže

$$x w + b \geq +1 \rightarrow \text{kategorie} = \text{správně}$$

$$x w + b \geq -1 \rightarrow \text{kategorie} = \text{nesprávně}$$

při čemž  $x$  je případ reprezentovaný jako vektor  $n$  komponent. Zde je  $w$  pomocný (Support) vektor, kolmý k nadrovině a odpovídá případům, které jsou mimo nebo nad limitem jejich kategorie (viz Graf 3).



Graf 3. Schema podpůrných vektorů.

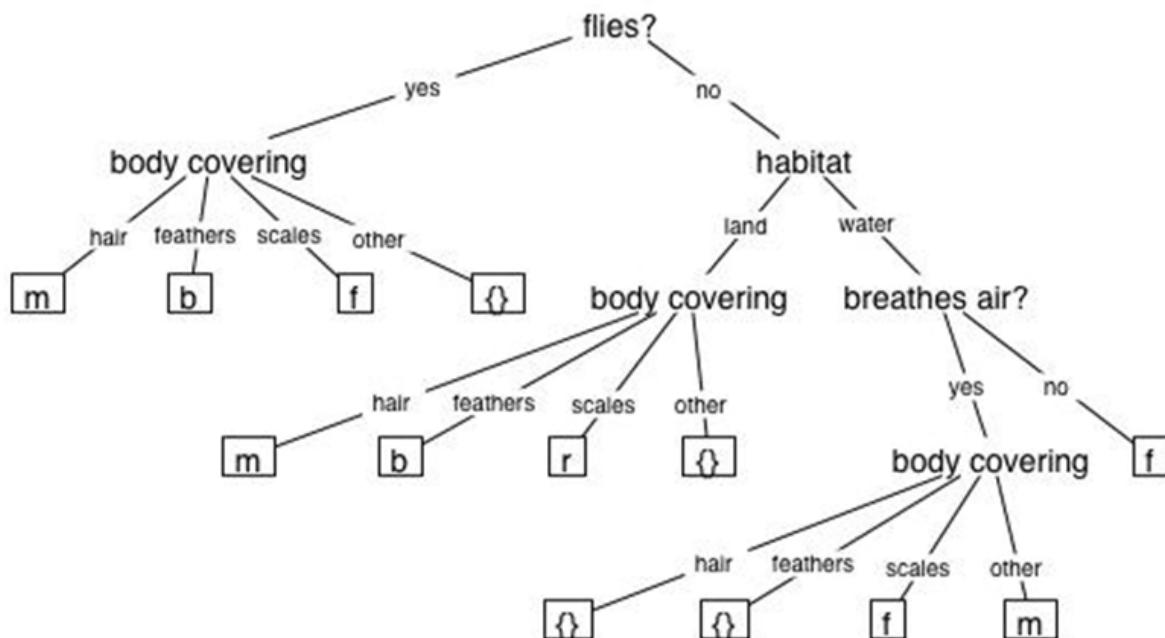
Pomocný vektor také svým modulem definuje rozpětí mezi nadrovinou a prvním pozitivním a prvním negativním případem (platí to pro prahy +1, -1). Pro každou kategorii se algoritmus snaží najít  $w$  s maximálním rozpětím, resp. se snaží najít nadrovinu, která má největší odstup od příkladů z trénovací množiny. Pro klasifikaci nezpracovaného případu se jednoduše použije výše uvedeného algoritmu. Tato jednoduchá realizace metody je jednou z metod používaných v našich pokusech, existují však další mnohem sofistikovanější přístupy a techniky.

### 3.4 K-nejbližší soused (K-Nearest Neighbour)

KNN (K-Nearest Neighbour) je algoritmus založený na paměti [15], jehož základní myšlenkou je, že zhodnocené případy nám analogií mohou pomoci řešit současný případ. Posuzuje každý případ jako vektor o  $n$  komponentách, přičemž  $n$  je počet atributů nebo charakteristik. Metoda nepotřebuje učící fázi. K předpovědi třídy řešeného případu porovná algoritmus řešeného případu se všemi případy trénovací množiny nebo s pamětí a vypočte vzdálenost mezi nimi. Potom je většinová třída pro  $K$  nejpodobnějších trénovacích případů předpovědí pro řešený případ. Vzdálenost použitá v případech je eukleidovská vzdálenost mezi vektory. Literatura uvádí řadu dalších možností.

### 3.5 ID3 a C4.5 rozhodovací stromy

Model vytvořený tímto algoritmem je strom [16], kde každý uzel odpovídá jednomu atributu a každá hrana odpovídá možné hodnotě uzlového atributu. Učící algoritmus vytváří strom z trénovacích dat. Výběr atributu, který bude tvořit uzel, vzniká v každém momentě výpočtem entropie dat po výběru uzlu. Pro každý atribut se počítá entropie zbývajících dat bez atributu podle různých hodnot uzlového atributu. Pro uzel je vybrán atribut, který vykazuje minimální entropii. Proces se opakuje do té doby, dokud nezbývá žádný atribut nebo je počet zbývajících případů pod uzlem menší než určitý limit. Příklad můžeme vidět na Grafu 4.



Graf 4. Příklad rozhodovacího stromu.

V příkladu jsou 4 atributy: *hmyz* (*flies*), *povrch těla* (*body covering*), *místo výskytu* (*habitat*), *dýchací průduchy* (*breathes air*), a 4 možné kategorie, *m*, *b*, *f* a *r*. První atribut je zde hmyz, protože je jediný, který rozděluje data na této úrovni s minimální entropií, atd. Pro klasifikaci řešeného případu stačí procházet strom směrem dolů a poslední listový atribut je předpovídaná kategorie. Na cestu od počátečního uzlu ke konečnému listovému atributu je možno nahlížet jako na pravidla, kde je pravidlo vytvořeno AND funkcí termínů (node=arc).

C4.5 je rozšíření ID3. Metoda umožňuje pracovat se spojitymi číselnými atributy, s chybějícími hodnotami a prořezávat strom při velkém množství dat. V experimentech použitý J48 strom je implementací C4.5.

### 3.6 RIDOR učení pravidel

RIDOR je název pro učení RIpple-DOwn pravidel [17]. Vytváří nejprve vstupní pravidlo a potom odchylky od vstupního pravidla s nejmenší (váženou) odchylkou a tímto klasifikuje trénovací data. Vznikají tak „nejlepší“ odchylky pro každou kategorii a proces se opakuje až do vyčištění souboru. Odchylky se jakoby řadí stromovitě a poslední list má pouze vstupní pravidla, ale žádné výjimky. Odchylky tvoří soubor pravidel, které předpovídají jiné třídy než jsou třídy ve vstupních pravidlech. Pro nalezení našich odchylek bylo užito algoritmu IREP. Tento algoritmus vytváří pravidla postupným přidáváním vždy jen jednoho výrazu do podmínek v každém kroku, takže počet omylů je minimalizován. Výraz pravidla podmínky je např. (*atribut {=, ≠, ≤, ≥} hodnota*).

## 4. Hodnocení

Hodnotící procesy a měření jsou stejné pro všechny pokusy. Z části dat byl vytvořen trénovací soubor, ostatní data byla testovaný soubor. Modely se učí z trénovacího souboru a snaží se předpovídat kategorie případů v testovaném souboru. Protože je známa kategorie dat testovaných případů, lze omezit předpovědi. Pro každou kategorii byly počítány tři různé typické míry: přesnost, úplnost a F-měření [18]. Přesnost je procento předpovědí v jedné kategorii, které jsou správné. Přesnost vyjadřuje Rovnice 1.

$$\text{Přesnost (kategorie}_i\text{)} = \frac{\text{počet správných předpovědí v kategorii}_i}{\text{celkový počet předpovědí v kategorii}_i} \quad (1)$$

Úplnost je procento všech případů v testovaném souboru v dané kategorii, které byly předpovědeny správně. Výpočet je uveden v Rovnici 2.

$$\text{Úplnost (kategorie}_i\text{)} = \frac{\text{počet předpovědí v kategorii}_i}{\text{celkový počet případů v kategorii}_i} \quad (2)$$

F-měření je kombinace výše uvedených výpočtů. Vyjadřuje jakoby průnik mezi případy, přesnost a úplnost standardizuje jejich součtem. Rovnice 3 ukazuje vyjádření F-měření.

$$F\text{-měření} = \frac{2 * \text{Přesnost} * \text{Úplnost}}{\text{Přesnost} + \text{Úplnost}} \quad (3)$$

Tyto tři míry byly počítány pro každou kategorii testovaného souboru. Jak bylo řečeno dříve, sesbíraná data je třeba rozdělit na soubor trénovací a testovaný. Společná cesta pro jejich vytvoření je křízová kontrola (cross-validation). Sesbíraná data jsou rozdělena na  $n$  stejně velkých částí. Každá  $n-1$  část kombinace je trénovací, zbytek je testovaný, algoritmus je  $n-1$  krát opakován a konečný výsledek je průměr z těchto  $n-1$  pochodů. Pro všechny níže popsané

pokusy, má  $n$  hodnotu 3, takže je vždycky 66 % trénovacích a 33 % testovaných dat a každý algoritmus se třikrát opakuje. Obvykle je hodnota  $n$  vyšší než 3, typicky to bývá 10, ale v tomto případě jsme měli v některých kategoriích velmi málo případů a vyšší hodnota  $n$  by mohla vytvořit soubor, který by nebyl reprezentativní pro danou kategorii, což je nežádoucí.

## 5. Pokusy

Byly provedeny dva druhy pokusů. První byl určený pro zjištění asociací mezi atributy, kde byla ukazatelem síly asociace charakteristika klasifikace. Druhý typ pokusů testoval předpověď budoucích příhod.

Je třeba poznamenat, že pozorování v souboru dat s chybějícími hodnotami nebyly vyloučeny z hodnocení ani nebyly dopočítávány, protože implementovaný učící algoritmus dokáže s chybějícími daty pracovat. Tyto implementace jsou vloženy do WEKA prostředí [19], které dokáže provést výše uvedené pokusy i s neúplnými daty.

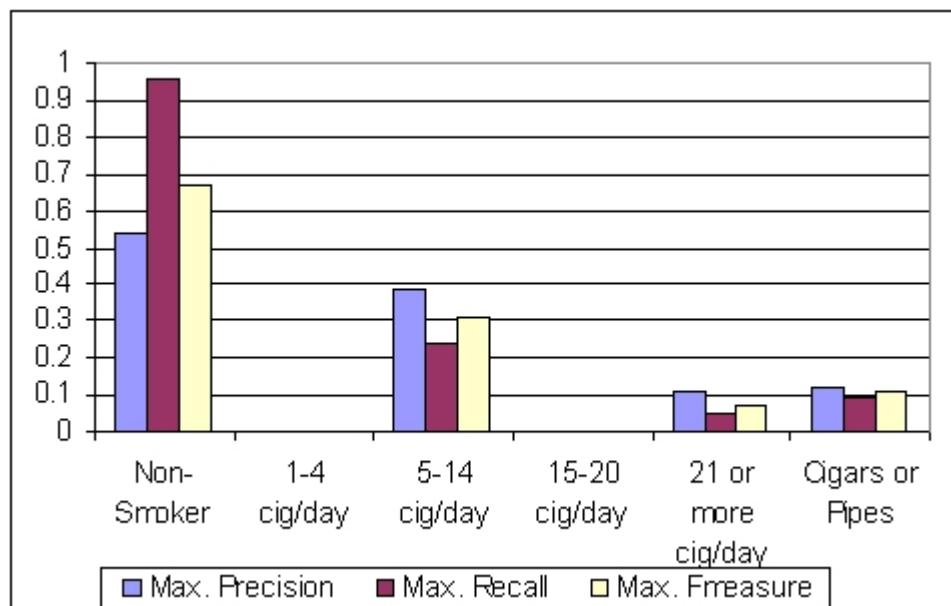
### 5.1 Nalezené odpovědi

První pokusy se vztahují k analytickým otázkám, které byly připraveny pro Výzvu k odhalování na konferenci ECML/OKDD 2004, konkrétně k těm, které se týkají souboru Vstup (Entry). Úlohy měly nalézt vztahy ve třech různých skupinách pacientů: normální skupina, riziková skupina a patologická skupina. Tyto skupiny odpovídaly riziku aterosklerózy – viz výše a jsou nazývány úrovní skupin. Konkrétně jsme hledali vybrané vztahy mezi sociálními faktory a tělesnou aktivitou, spotřebou alkoholu, kouřením, hmotnostním indexem, krevním tlakem a HDL cholesterolom. Další úloha měla zjistit vztah mezi tělesnou aktivitou a ostatními faktory a mezi spotřebou alkoholu a ostatními faktory. Data každé skupiny byla použita v algoritmu strojového učení se snahou předpovědět hodnotu každého faktoru v jedné skupině podle hodnot tohoto faktoru v druhé skupině s ohledem na možné hodnoty uvažované kategorie. Např. čtyři sociální faktory byly vloženy do algoritmu jako trénovací faktory, aby bylo možno předpovídat hodnotu každého atributu tělesné aktivity apod. pro ostatní skupiny faktorů. Pro každý vztah byly počítány maximální hodnoty z výsledků různých algoritmů, aby bylo možno provést porovnání mezi úrovněmi skupin. Pokud byla předpověď přesná, lze říci, že jde o silný vztah mezi faktory trénovacími a faktorem, jehož hodnota byla předpovídána, a to se stupněm síly odpovídajícím přesnosti předpovědi. Lze také porovnat míru předpovědi mezi faktory a úrovní skupiny ke stanovení, které vztahy jsou silnější než jiné.

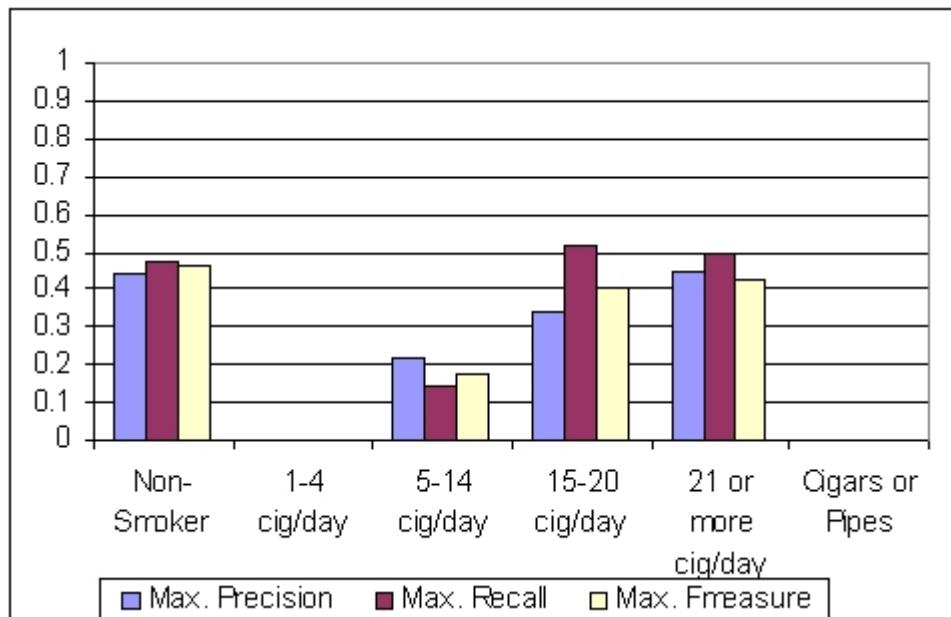
Vzhledem k omezené délce sdělení jsou prezentovány jen některé typické výsledky. Na Grafu 5 jsou uvedeny výsledky předpovědi, a to míry maximální přesnosti, úplnosti a F-měření pro sociální atribut „Kouření (Smoking)“ proti ostatním sociálním faktorům v každé úrovni skupin, a to a) Normální, b) Patologické a c) Rizikové, a proti tělesné aktivitě ve skupinách d) Normální, e) Patologické a f) Rizikové. Je vidět, že pro Normální skupinu jsou nejlepší předpovědi pro nekuřáky, a to jak v sociálních faktorech tak v tělesné aktivitě, zatímco pro ostatní stupně „Kouření“ byly výsledky nevýznamné. Zdá se, že pro vztah mezi sociálními faktory a kouřením je vztah poněkud silnější než je vztah mezi tělesnou aktivitou a kouřením, protože jsou lepší výsledky ve všech stupních atributu kouření. V Patologické a Rizikové skupině je vztah mezi trénovacími faktory a nekuřáctvím silnější než je vztah mezi kouřením a tělesnou aktivitou – zvláště silný byl tento vztah v Patologické skupině. V této skupině byly

Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze

osoby, které kouří 15 nebo více cigaret za den lépe předpověděny než v Normální skupině, ale nekuřáci byli detekováni mnohem hůř než v Normální skupině.

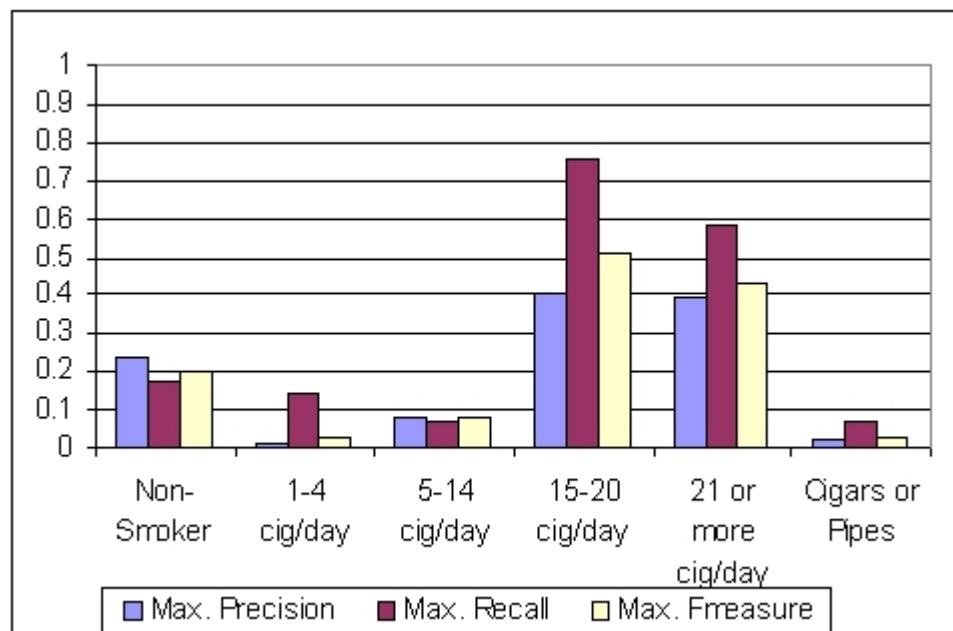


a)

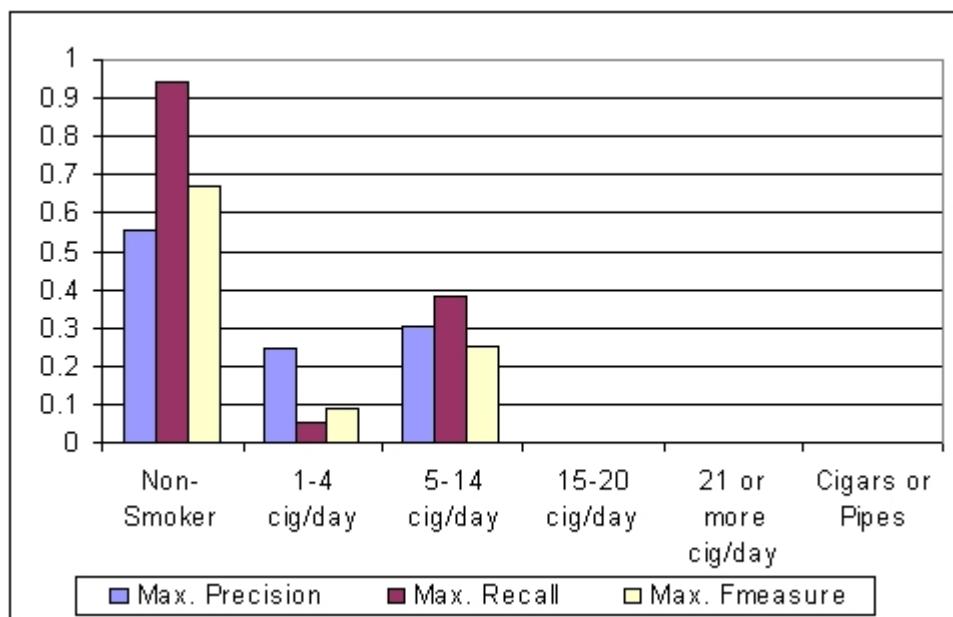


b)

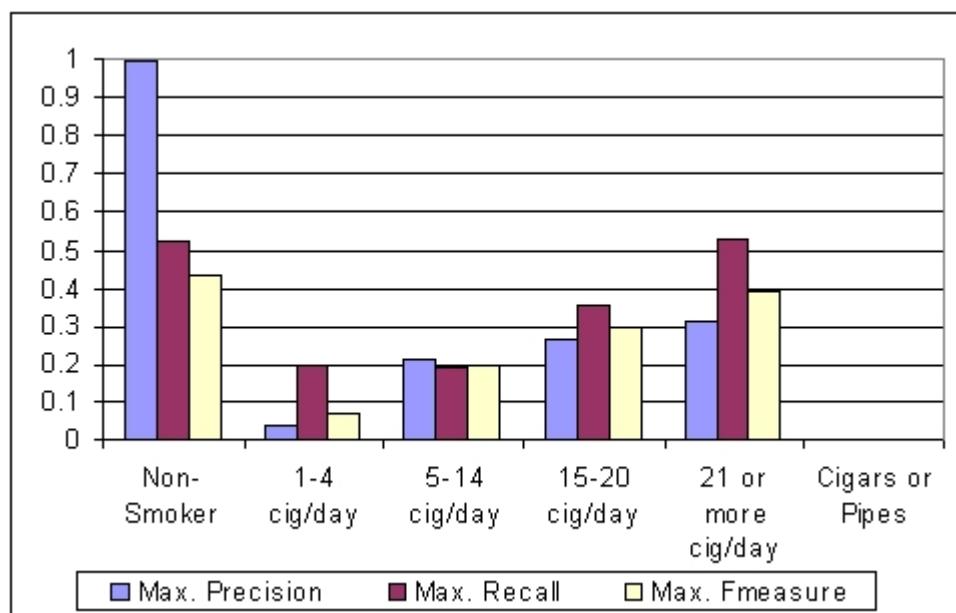
Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze



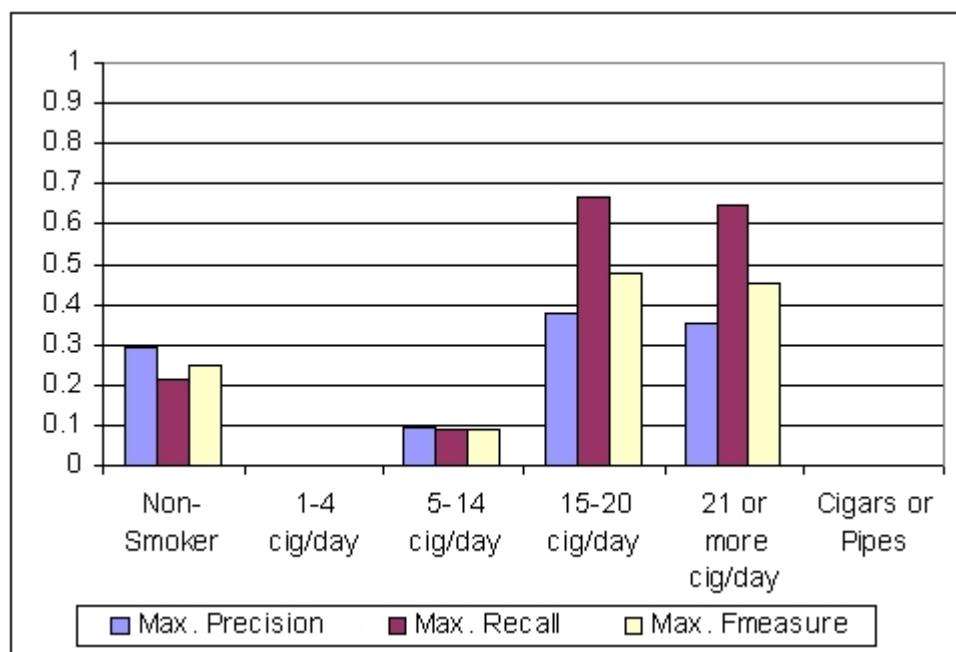
c)



d)



e)



f)

Graf 5. Hodnoty maximální přesnosti, maximální úplnosti a maximálního F-měření ve všech algoritmech pro předpověď atributu „Kouření (Smoking)“ z ostatních sociálních faktorů v a) Normální skupině, b) Patologické skupině, c) Rizikové skupině a z faktorů tělesné aktivity v d) Normální skupině, e) Patologické skupině, f) Rizikové skupině.

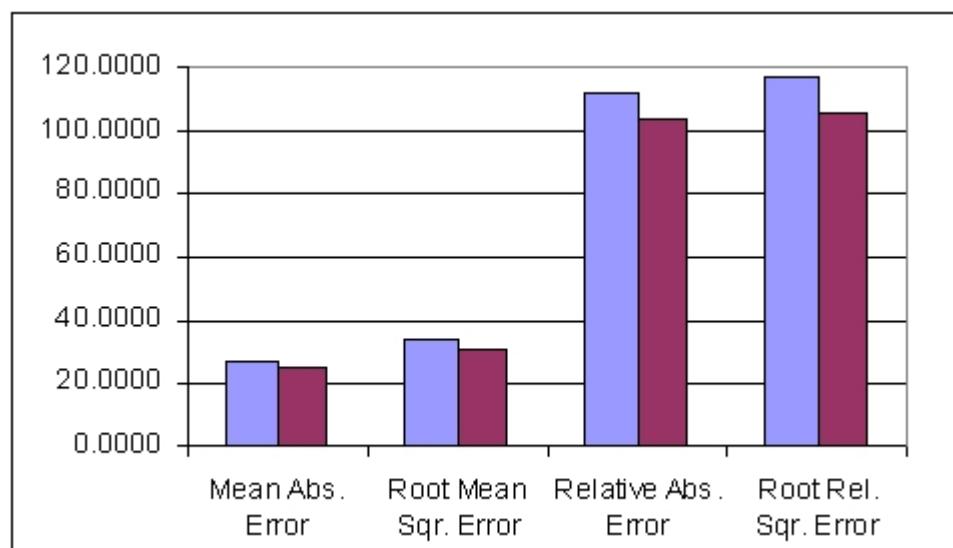
Vysvětlivky:

Max. Precision – maximální přesnost, Max. Recall – maximální úplnost, Max. Fmeasure – maximální F-měření

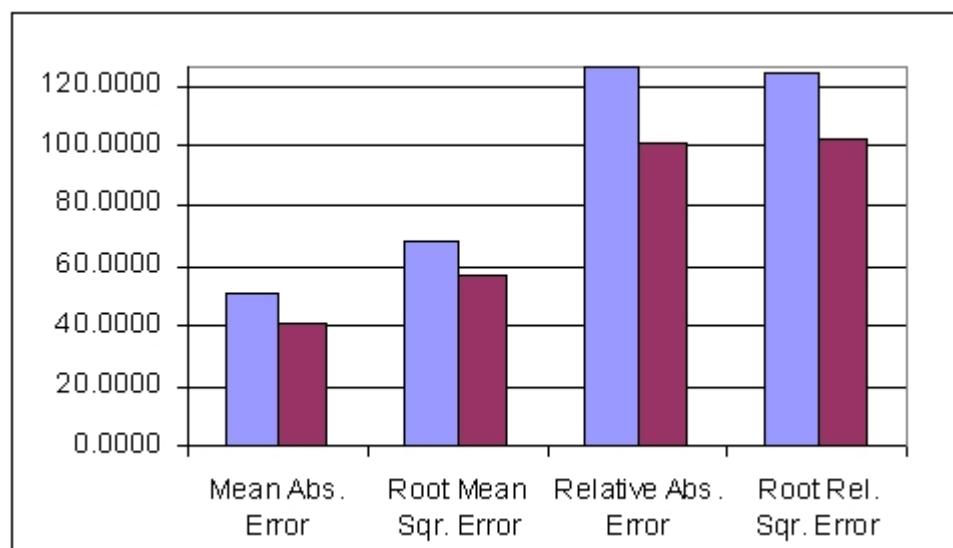
Non-Smoker – nekuřáci, 1-4 cig/day a další – kuřáci 1-4 cigaret/den a další, Cigars or Pipes – kuřáci doutníků nebo dýmků.

Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze

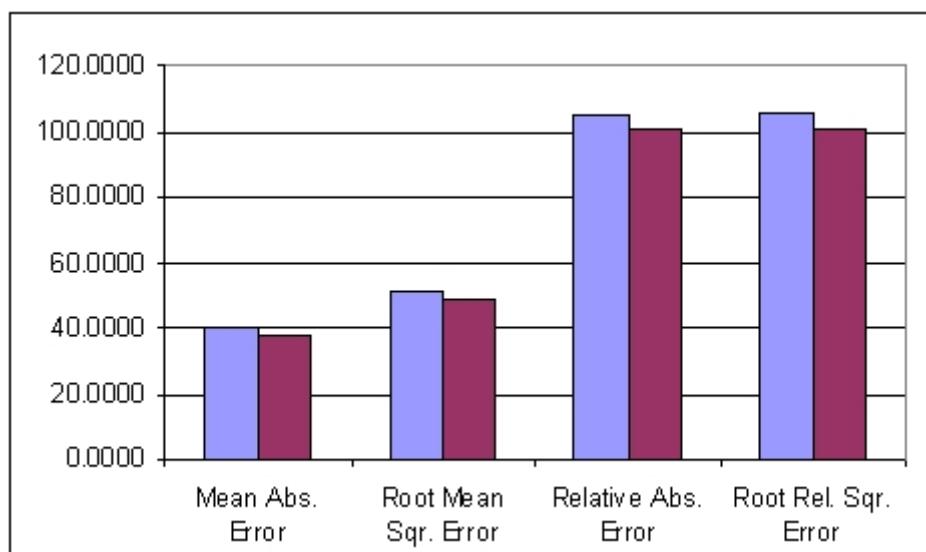
Podívejme se na jiný charakteristický příklad. Graf 6 ukazuje výsledky předpovědi hladiny cholesterolu ze sociálních faktorů – a), b), a c) a z faktorů tělesné aktivity, d), e) a f) pro každou úroveň skupin. V tomto případě jsou výsledky předpovědi velmi podobné ve vztahu mezi sociálními faktory a cholesterolom a mezi faktory tělesné aktivity a cholesterolom ve všech úrovních skupin, takže můžeme říci, že síla těchto vztahů je také podobná. Nicméně se liší mezi úrovněmi skupin. V Normální skupině je průměrná absolutní chyba předpovědi okolo 24, zatím co v Patologické a Rizikové skupině je okolo 50, resp. 40. Lze učinit závěr, že je snadnější předpovědět hladinu cholesterolu z obou skupin faktorů – sociálních i faktorů tělesné aktivity pro osoby z Normální skupiny než pro osoby z ostatních dvou skupin. Ukazuje to na silný vztah mezi trénovacími faktory a hladinou cholesterolu v posledně jmenované skupině.



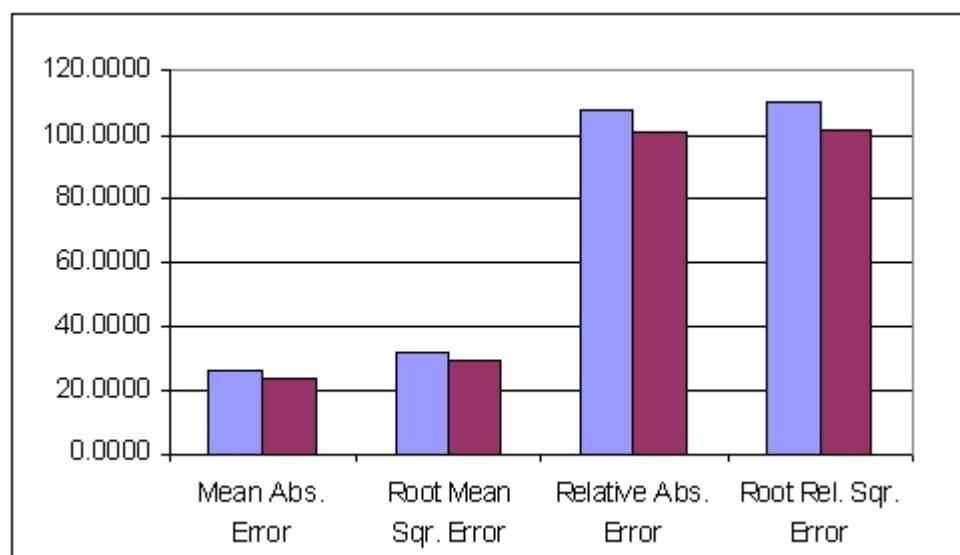
a)



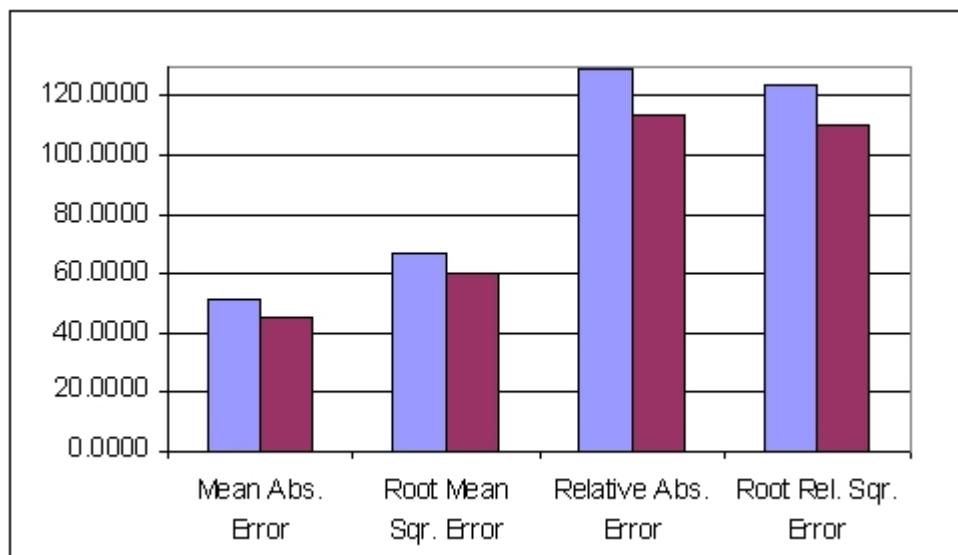
b)



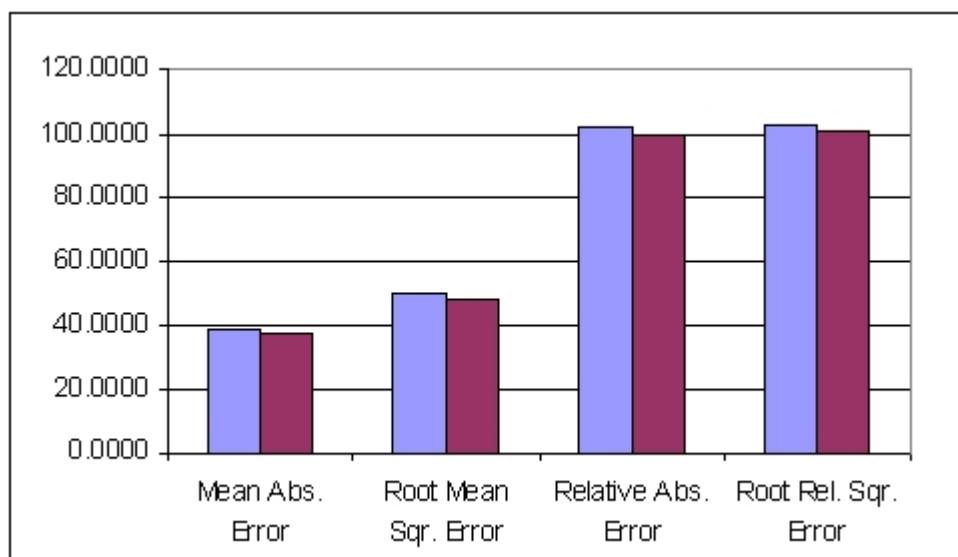
c)



d)



e)



f)

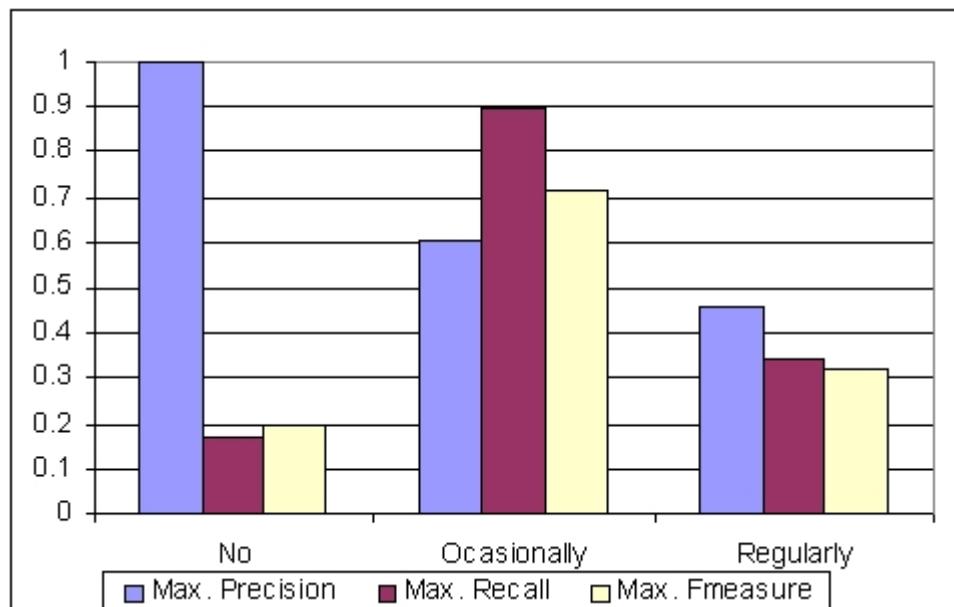
Graf 6. Hodnoty průměrné absolutní chyby, odmocniny střední kvadratické chyby, podílu absolutních chyb a odmocniny relativní kvadratické chyby ze všech algoritmů předpovědi hladiny cholesterolu ze sociálních faktorů v a) Normální skupině, b) Patologické skupině a c) Rizikové skupině a z faktorů tělesné aktivity v d) Normální skupině, e) Patologické skupině a f) Rizikové skupině.

Vysvětlivky:

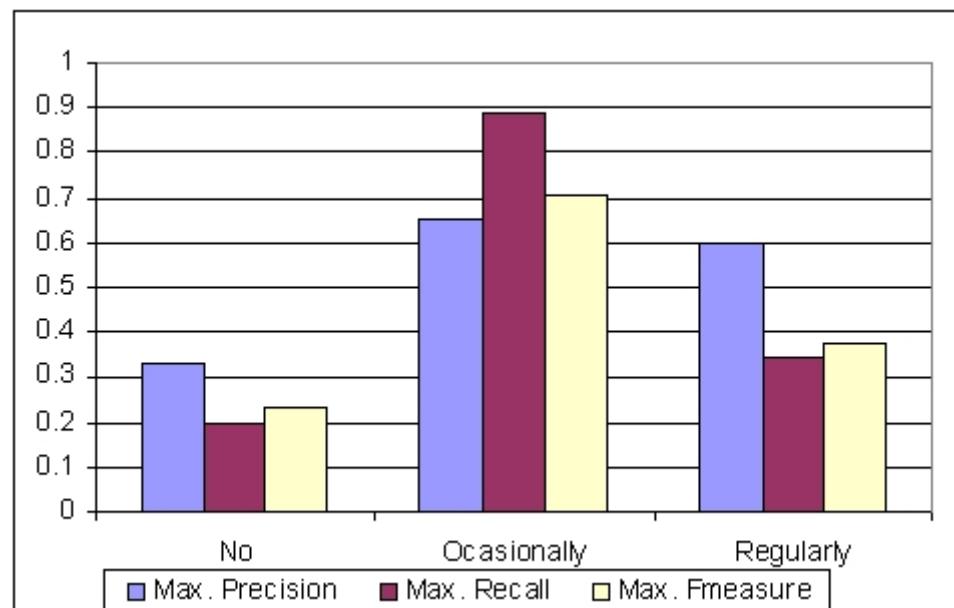
Mean Abs.Error – průměr absolutní chyby, Root Mean Sqr. Error – odmocnina střední kvadratické chyby, Relative Abs. Error – podíl absolutních chyb, Root Rel. Sqr. Error – odmocnina relativní kvadratické chyby.

Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze

A konečně, Graf 7 ukazuje výsledky pro předpověď hodnot příjmu alkoholu, zvlášť ze sociálních faktorů a z faktorů tělesné aktivity jako trénovacích, pro každou úroveň, podobně jako bylo uvedeno výše.

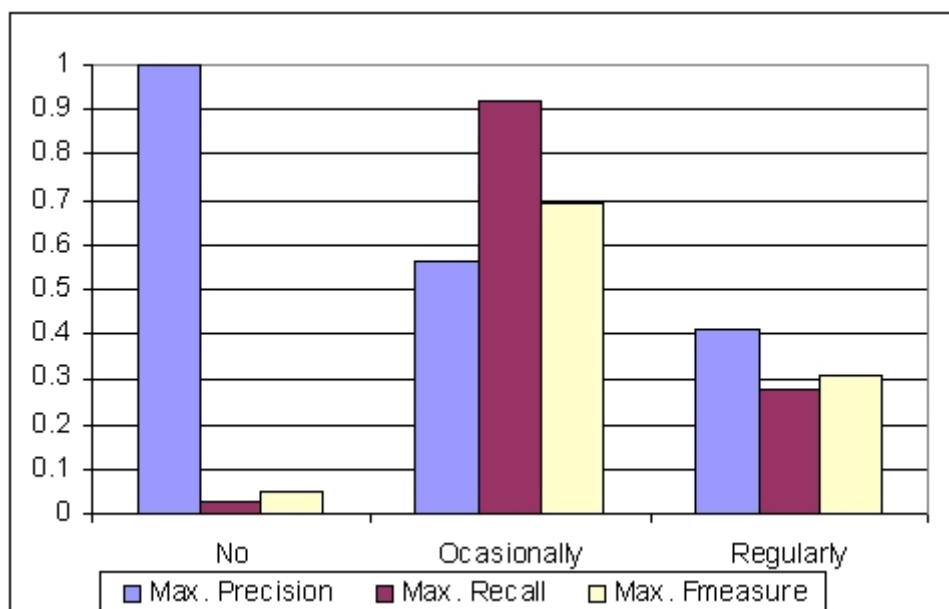


a)

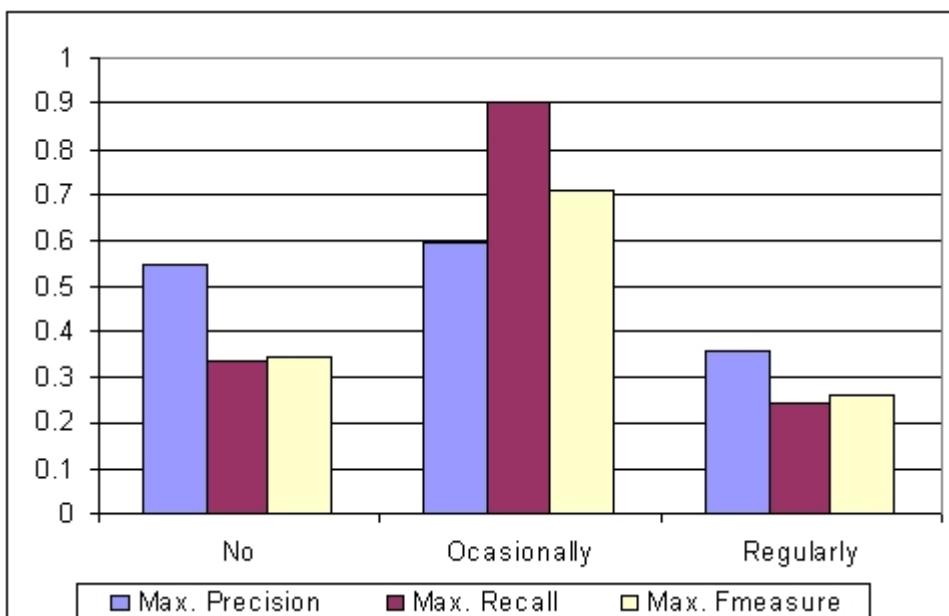


b)

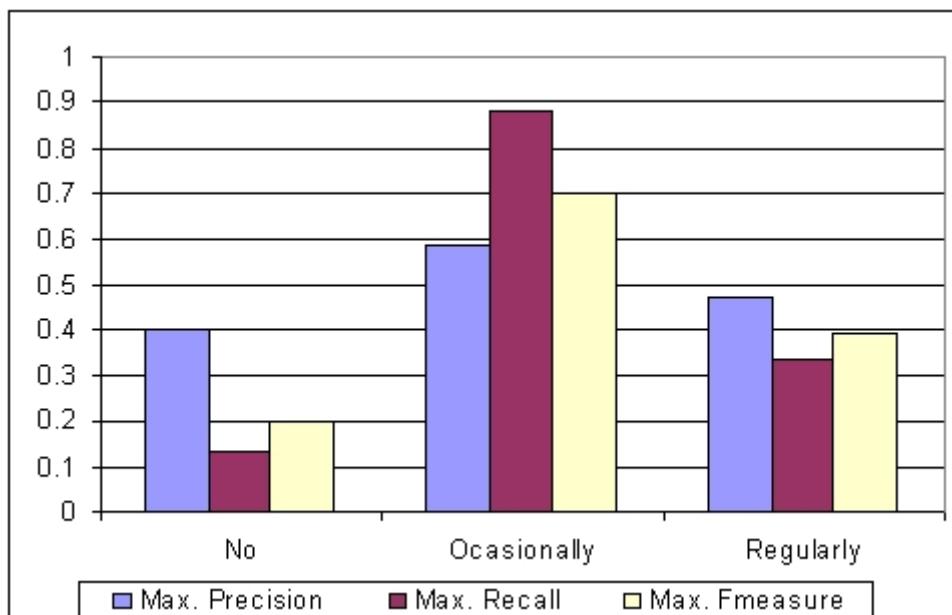
Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze



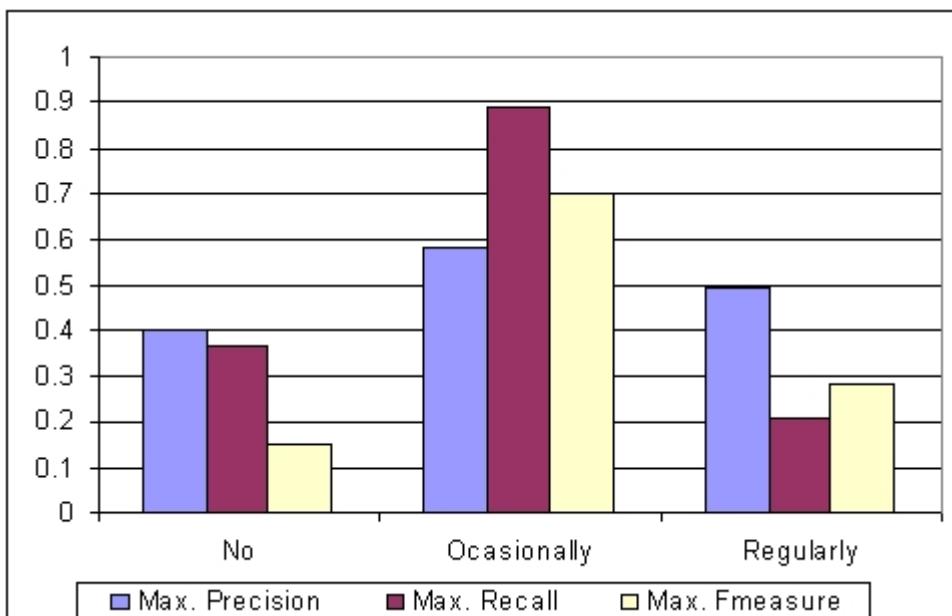
c)



d)



e)



f)

Graf 7. Hodnoty maximální přesnosti, maximální úplnosti a F-měření ze všech algoritmů pro předpověď atributu „Alkohol“ pouze ze sociálních faktorů v a) Normální skupině, b) Patologické skupině, c) Rizikové skupině a pouze z faktorů tělesné aktivity v d) Normální skupině, e) Patologické skupině, f) Rizikové skupině.

Vysvětlivky:

Max. Precision – maximální přesnost, Max Recall – maximální úplnost, Max Fmeasure – maximální F-měření.

### Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze

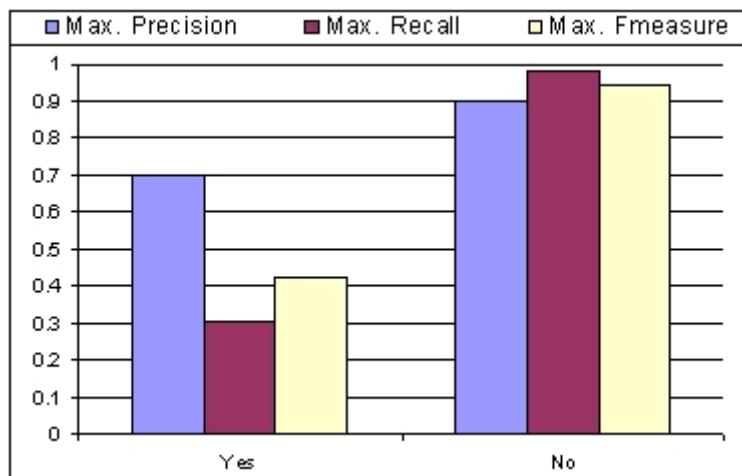
Výsledky v Grafu 7 ukazují, že je ve všech úrovních skupin jasný vztah mezi trénovacími faktory a osobami, které pijí alkohol příležitostně. Osoby, které pijí alkohol pravidelně, jsou mnohem hůře detekovatelní a hůře se z trénovacích faktorů předpovídají, takže tento vztah je nevýrazný, snad poněkud silnější v Patologické skupině. Totéž lze říci o vztahu osob, které alkohol nepijí a jejich tělesné aktivitě. Nicméně významně vyšší je přesnost předpovědi ze sociálních faktorů v Normální a v Rizikové skupině. Osoby, které nepijí alkohol, jsou přesně identifikovány ze sociálních faktorů v Rizikové skupině, což ukazuje na významný vztah mezi danými atributy.

Trénovací skupinu faktorů tvořily všechny atributy dohromady. Z lékařského pohledu je také zajímavé tyto faktory oddělit a vytvořit z nich podsoubor. Takže se např. předpovídala tělesná aktivita v zaměstnání ze všech možných kombinací sociálních faktorů. Výsledky ukázaly, že pro Normální skupinu a Rizikovou skupinu dává mnohem lepší předpovědní výsledek samotný faktor „Vzdělání (Education)“, než jakákoli jiná kombinace sociálních faktorů. V Patologické skupině je to podobné, ale rozdíl není tak velký jako v ostatních skupinách, nejlepší kombinace je v této skupině „Věk + Vzdělání (Age + Education)“.

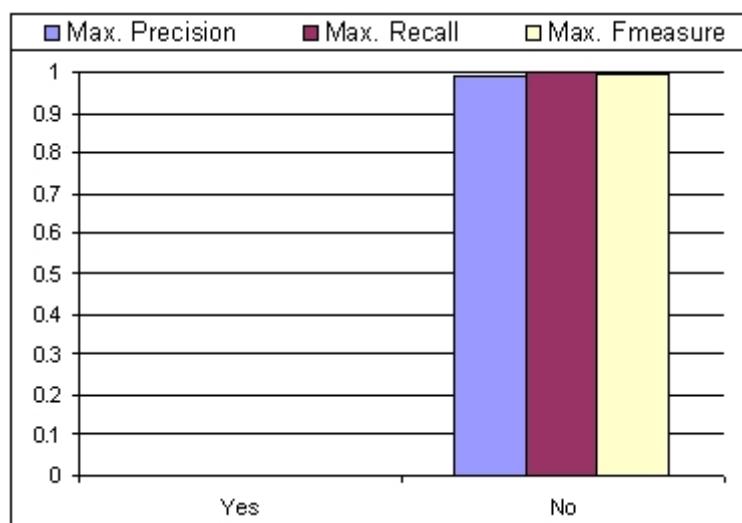
## 5.2 Předpověď budoucích příhod

Hlavní cíl následujících pokusů bylo testovat přesnost předpovědních algoritmů. Bylo užito souboru Vstup (Entry), ale i Kontroly (Control). Nejprve byli vybráni pacienti, kteří měli kontrolní záznam v souboru Control po deseti letech od vstupu do studie. Poté byl proveden pokus předpovědět z atributů Vstup (Entry), zda u nich vznikne nějaká příhoda. Jako příhody byly vybrány systolicko-diastolická hypertenze, systolická hypertenze, diastolická hypertenze, hypercholesterolemie a hypertriglyceridemie. Možné hodnoty těchto atributů byly „správně“ nebo „nesprávně“. Stejně tak byly zhodnoceny záznamy po dvaceti letech. Výsledky ukázaly, že nejlepším algoritmem byl vícevrstevní perceptron, který dosáhl téměř 85% přesnosti a 65% úplnosti v detekci všech poruch. Zatímco riziko vzniku hypertenze bylo v Rizikové skupině 0, protože někteří pacienti této skupiny měli hypertenzi do začátku studie, je z lékařského hlediska mnohem zajímavější provést tyto pokusy jen pro Normální skupinu. Takže pro tuto skupinu byl proveden stejný proces pro deset a pro dvacet let. Výsledky pro různé uvažované poruchy ukazuje Graf 8, a to a) až e) pro desetileté předpovědi a f) až j) pro dvacetileté předpovědi. Pro každou příhodu je uvedena maximální hodnota ze všech použitých algoritmů. V tomto případě výsledky ukázaly, že není jeden nejlepší algoritmus. Podle předpovídání příhody a pro některé kategorie je některý algoritmus lepší než ostatní (ukazované maximální hodnoty odpovídají různým algoritmům), takže je zajímavé použít všechny algoritmy a rozhodnout podle výsledků všech algoritmů. Nutno upozornit na to, že přesnost předpovědi je mnohem vyšší pokud jsou pro vstup dána data všech tří skupin dohromady s tím, že první zájem je na Normální skupině.

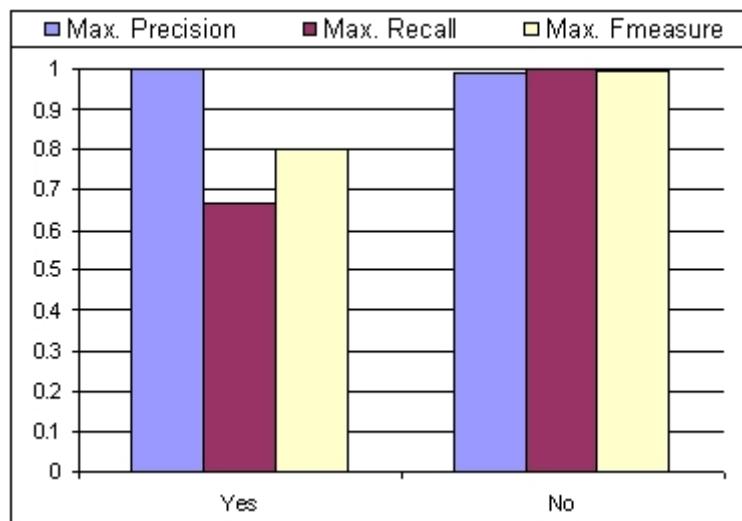
Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze



a)



b)

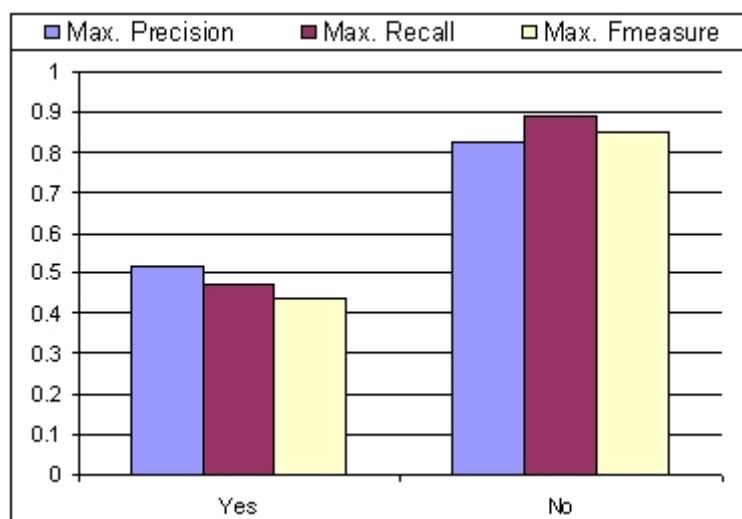


c)

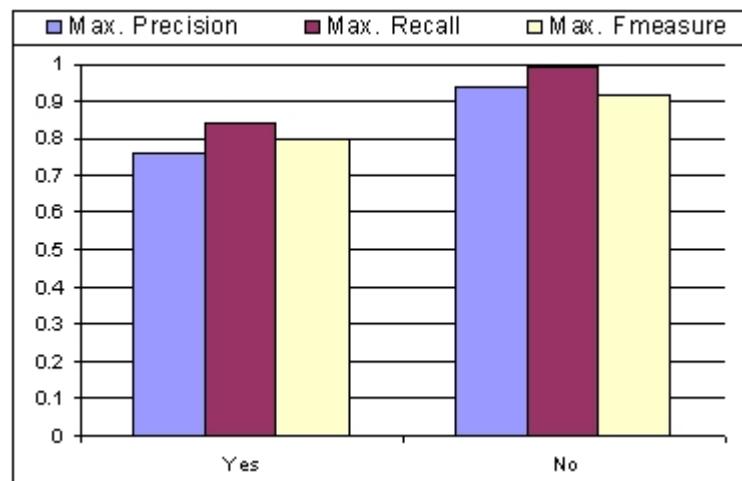
Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze



d)



e)

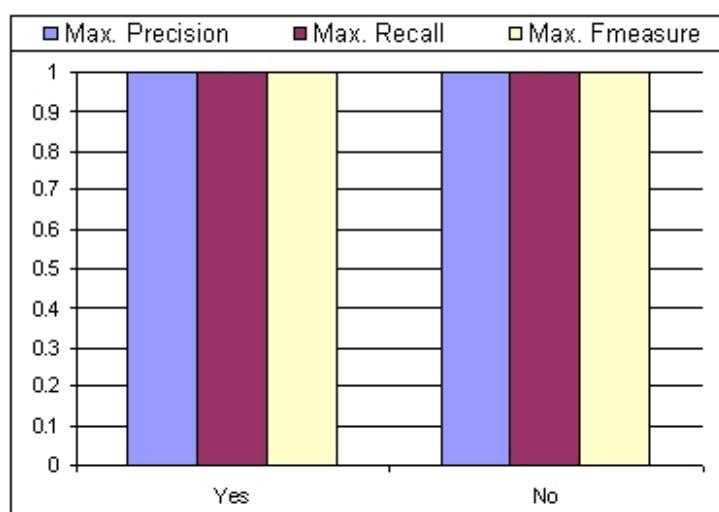


f)

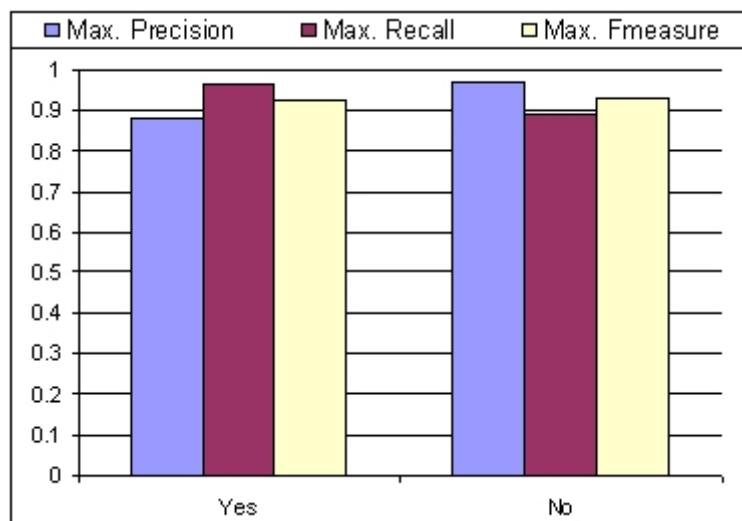
Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze



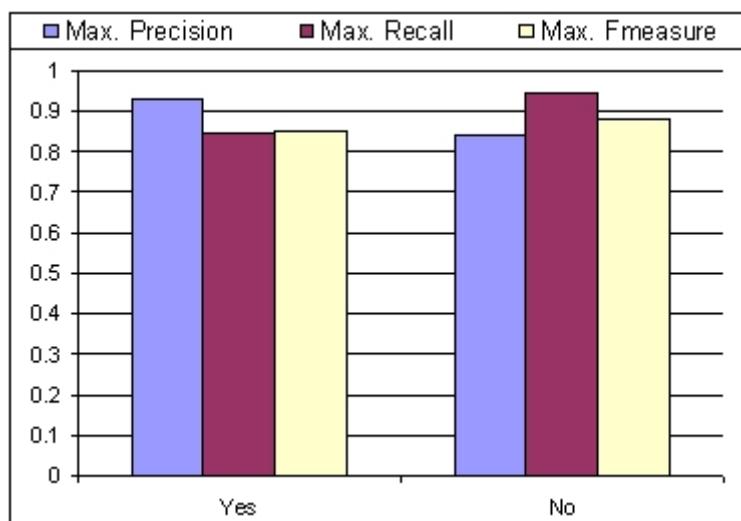
g)



h)



i)



j)

Graf 8. Hodnoty maximální přesnosti, úplnosti a F-měření pro předpověď a) systolicko-diastolické hypertenze, b) systolické hypertenze, c) diastolické hypertenze, d) hypercholesterolemie a e) hypertriglyceridemie v deseti letech, a f) systolicko-diastolické hypertenze, g) systolické hypertenze, h) diastolické hypertenze, i) hypercholesterolemie a j) hypertriglyceridemie ve dvaceti letech.

Vysvětlivky:

Max. Precision – maximální přesnost, Max Recall – maximální úplnost, Max Fmeasure – maximální F-měření.

Hodnoty v Grafu 8 ukazují, že předpověď příhody je přesněji odvozena po 20 letech, ale velmi špatně předpověděna po 10 letech s výjimkou diastolické hypertenze. Nepřítomnost příhody je stejně dobře předpověděna pro deset jako pro dvacet let. Ze všech příhod vychází nejlepší předpověď pro diastolickou hypertenzi, kde je hodnota přesnosti předpovědi téměř 100 % pro její přítomnost i pro její nepřítomnost. Nejhorší byla předpověď pro systolickou hypertenzi – přítomnost v deseti letech nebyla detekovatelná.

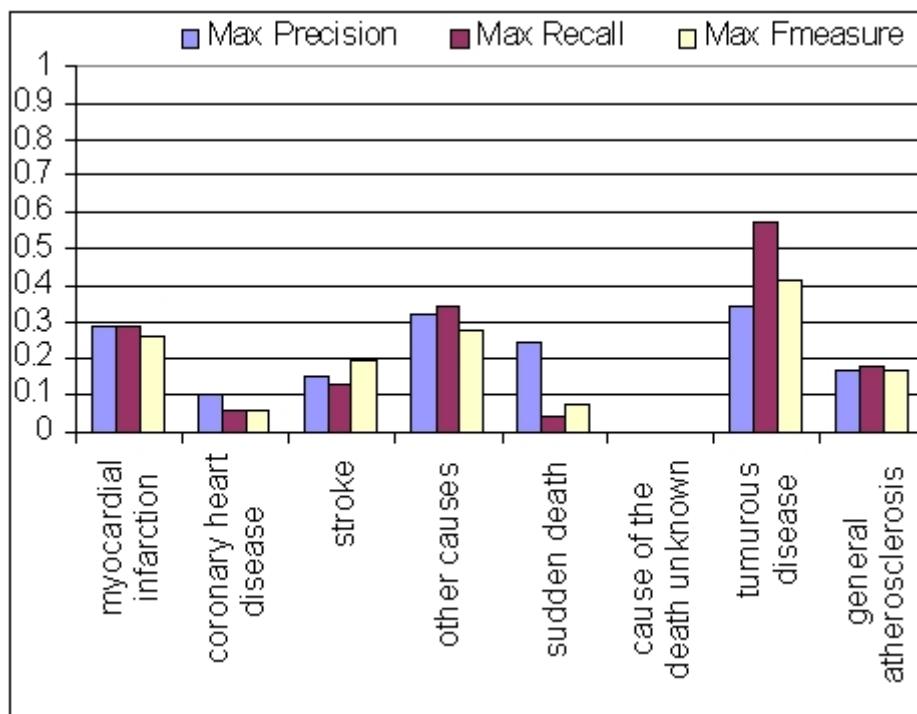
Pro předpověď některých dalších příhod, např. anginy pectoris, srdečního infarktu, mozkových příhod a dalších nebylo možné pro malý počet pozorování těchto příhod, takže výsledky nejsou relevantní.

### 5.3 Předpověď příčiny úmrtí

Tento pokus je podobný jako předcházející, byla zde ale předpovídána příčina úmrtí místo onemocnění nebo poruch. Bylo užito souboru Úmrtí (Death). Algoritmy byly pro zemřelé osoby trénovány na jejich datech ze Vstupu (Entry). Pokusy byly provedeny samostatně pro tři úrovně skupin a poté pro všechna vstupní data všech tří úrovní skupin dohromady. Výsledky ukazuje Graf 10. V Normální skupině – Graf 10b) byla nejlépe předpověděnou příčinou nádorová onemocnění a ostatní příčiny úmrtí. V Rizikové skupině – Graf 10 d), byla nejlepší předpověď pro ostatní příčiny, ale též pro srdeční infarkt a ischemickou chorobu srdeční, které nebyly vůbec předpověděny v Normální skupině. V Patologické skupině – Graf 10 c), byly nejlépe předpověděny srdeční infarkt a nádorová onemocnění, ale předpověď

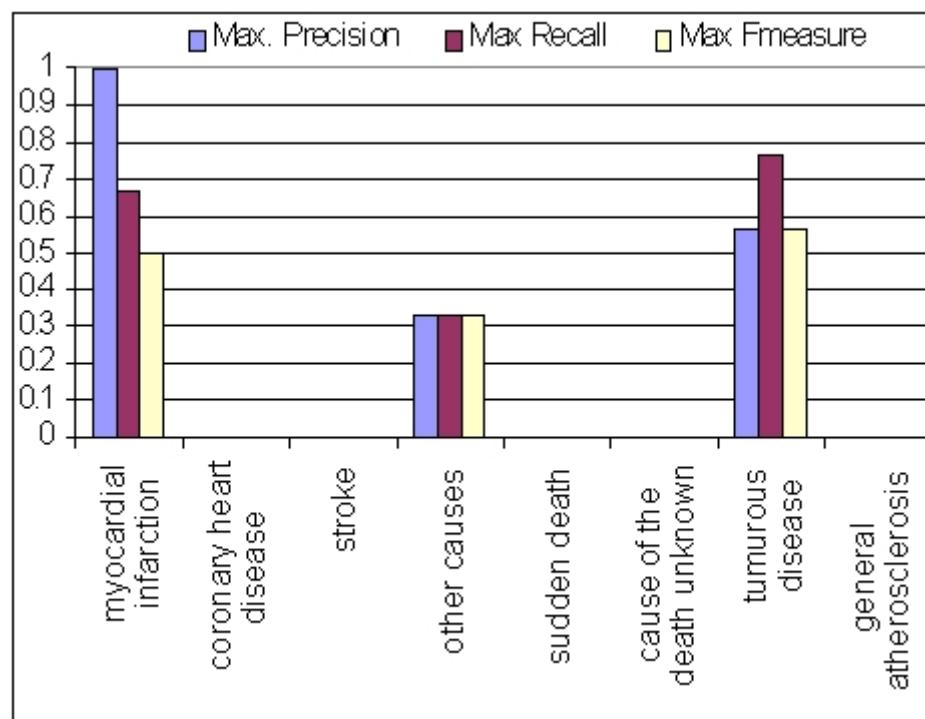
Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze

mozkové mrtvice a celkové aterosklerózy byla nedokonalá, dostali jsme mnohem horší výsledky pro tyto příhody než v ostatních skupinách. Souhrnně – Graf 10a), předpověď příčiny úmrtí byla velmi nedokonalá a to proto, že data ze souboru Vstup (Entry) neměla dostatek informací pro předpověď příčiny úmrtí a/nebo by možná bylo třeba více pozorování. Ale, kolik je to dostatek informací?

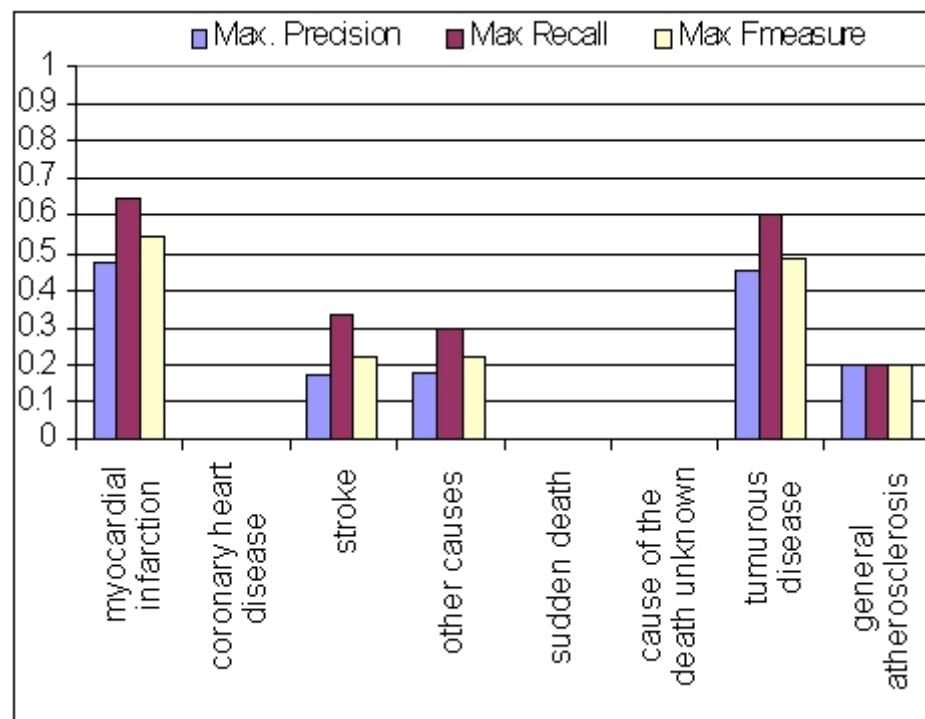


a)

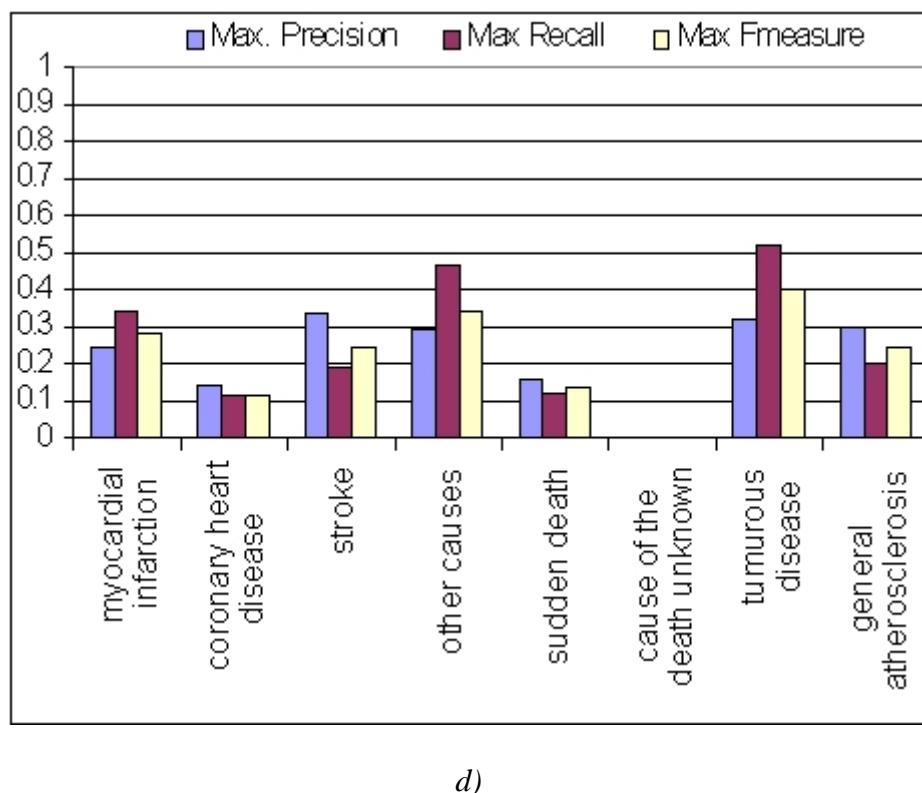
Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze



b)



c)



d)

Graf 9 Hodnoty maximální přesnosti, výpovědní hodnoty a F-měření pro předpověď příčiny úmrtí, a) všechny skupiny dohromady, b) Normální skupina, c) Patologická skupina a d) Riziková skupina.

Vysvětlivky:

Myocardial infarction – srdeční infarkt, coronary heart disease – ischemická choroba srdeční, stroke – mozková mrtvice, other causes – jiné příčiny, sudden death – náhlá smrt, cause of the death unknown – neznámá příčina úmrtí, tumorous disease – nádorové onemocnění, general atherosclerosis – celková ateroskleróza, Ostatní - viz vysvětlivky u Grafu 8.

## 6. Závěry

Různé typy algoritmů strojového učení byly použity k vyhledávání znalostí z lékařských dat, a to dvojím způsobem: za prvé, metody byly použity k předpovědi hodnoty jednoho atributu z databáze pacienta, zatímco ostatní atributy vytvořily trénovací soubor. Záměrem bylo stanovit maximální přesnost mezi všemi algoritmy jako míru síly vztahu mezi trénovací částí a cílovým atributem. Toto měření se ukázalo užitečné i pro porovnání vztahů mezi atributy v různých skupinách pacientů.

Za druhé, učící techniky byly použity k předpovědi budoucích příhod. Výsledky ukázaly, že některé metody předpovídají některé příhody lépe než ostatní, takže je zajímavé použít všechny algoritmy najednou a zhodnotit spolehlivost výsledků podle známého trendu každé metody. Všechny testované metody poskytly lepší předpověď pro dvacet let než pro deset let sledování a pro některé příhody dosáhly výborných výsledků, takže by to mohly být metody vhodné pro podporu rozhodování. Algoritmy strojového učení byly použity také pro

Metody strojového učení pro vyhledávání znalostí v lékařských datech o ateroskleróze

předpověď příčin úmrtí, v tomto případě bylo dosaženo špatných výsledků, možná pro malé množství informací (vstupů) v tomto souboru dat.

Do budoucna by mohlo být zajímavé zjemnit nastavení parametrů algoritmů a testovat více technik. V úmyslu je také spojit všechny významné a použitelné metody z této práce a vytvořit expertní systém a výzkumně odvodit z výsledků systému srozumitelná pravidla.

## Poděkování

Výzkum byl částečně podpořen Výzkumným plánem Ústavu informatiky AV ČR AV0Z10300504 a Výzkumným plánem Španělské Rady pro Vědecký výzkum spolu s podporou Ústavu průmyslové automatizace „María Bueno“.

## Literatura

- [1] Mitchell, T.: Machine Learning. McGraw Hill, 1997.
- [2] Lavrać, N.: Selected Techniques for Data Mining in Medicine. Artificial Intelligence in Medicine, vol. 16 (1), pp. 3-23, 1999.
- [3] Aseervatham, S. and Osmani A.: Mining Short Sequential Patterns for Hepatitis Type Detection. ECML/PKDD Discovery Challenge, 2005.
- [4] Aubrecht, P., Kejkula, M., Kremen, P., Novakova, L., Rauch, J., Simunek, M., Stepankova, O.: Mining in Hepatitis Data by LISp-Miner and SumatraTT. ECML/PKDD Discovery Challenge, 2005.
- [5] Pizzi, L.C., Ribeiro, M.X., Vieira, M.T.P.: Analysis of Hepatitis Dataset using Multirelational Association Rules. ECML/PKDD Discovery Challenge, 2005.
- [6] Durand, N., Soulet, A.: Emerging Overlapping Clusters for Characterizing the Stage of Liver Fibrosis. ECML/PKDD Discovery Challenge, 2005.
- [7] Durand, N., Cleuziou, G., Soulet, A.: Discovery of Overlapping Clusters to Detect Atherosclerosis Risk Factors. ECML/PKDD Discovery Challenge, 2004.
- [8] Cios, K. J.: Medical data mining and Knowledge Discovery. Physica – Verlag, 2001.
- [9] Chen, H., Fuller, S. S., Friedman, C. and Hersh, W.: Medical Informatics: Knowledge Management and Data Mining in Biomedicine. Integrated Series in Information Systems (2), Springer Science and Business Media Inc., 2005.
- [10] Boudik F., Reissigová J., Hrach K., Tomecková M., Bultas J., Anger Z., Aschermann M., Zvarová J.: Primary Prevention of Coronary Artery Disease Among Middle Aged Men in Prague: Twenty-year Follow-up Results. Atherosclerosis. 2006 Jan;184(1):86-93.
- [11] Tomecková, M.: The Challenge on Atherosclerosis Data Viewed by the Experts. ECML/PKDD Discovery Challenge, 2004.
- [12] Rish, I.: An Empirical Study of the Naive Bayes Classifier. IJCAI-01 Workshop on Empirical Methods in AI, 2001.
- [13] Haykin, S.: Neural Networks: A comprehensive Foundation (2nd edition). Pearson Education, 1998.

- [14] Scholkopf, B., Smola, A. J., Müller, K.-R., Burges, C. J. C., and Vapnik, V.: Support Vector Methods in Learning and Feature Extraction. In Down, T., Frean, M., and Gallagher, M., editors. Proceedings of the Ninth Australian Congress on Neural Networks, Brisbane, Australia. University of Queensland, 1998.
- [15] Teknomo, K.: K-Nearest Neighbors Tutorial. <http://people.revoledu.com/kardi/tutorialKNN>, 2004.
- [16] Quinlan, J. R.: C4.5: Programs for Machine Learning. Morgan Kauffman, 1993.
- [17] Compton, P., Edwards, G., Kang, B., Malor, R., Menzies, T., Preston, P., Srinivasan, A. and Sammut, S.: Ripple Down Rules: Possibilities and Limitations. Boose, J.H. & Gaines, B.R., Ed. Proceedings of the Sixth AAAI Knowledge Acquisition for Knowledge-Based Systems Workshop. pp.6-1-6-20. Calgary, Canada, University of Calgary, 1991.
- [18] Van Rijsbergen, C. J.: Information Retrieval. Butterworths, London, 1979.
- [19] Witten, I. H. and Frank, E.: Data Mining: Practical Machine Learning Tools and Techniques. 2nd Edition, Morgan Kaufmann, San Francisco, 2005.

# Low-dimensional Multimodal Deformable Registration of MRI Brain Images in Stereotaxic Space

Daniel Schwarz<sup>1</sup>, Ivo Provazník<sup>2</sup>

1. Institute of Biostatistics and Analyses, Masaryk University, Czech Republic,

2. Department of Biomedical Engineering, Brno University of Technology, Czech Republic

Deformable image registration is a fundamental technique in computational neuroanatomy. An iterative multilevel block matching technique with the use of several recent inventions is proposed here. A symmetric multimodal similarity measure allows to register subject images to an arbitrary digital brain atlas. Smooth deformations produced by scattered data interpolation based on compactly supported radial basis functions suppress gross inter-subject differences and preserve the localized anatomical variability which may be further studied with selected automated morphometry methods. Four similarity measures are tested in an experiment with image data obtained from the Simulated Brain Database and a quantitative evaluation of the algorithm is presented.

**Keywords:** image processing, image registration, MRI images, computational neuroanatomy, radial basis functions

## Introduction

One of widely applied methods in computational neuroanatomy is a voxel-based morphometry (VBM), which has recently become a subject of discussion [1], [2]. It interrogates anatomical MRI scans on voxel by voxel basis, in order to demarcate regions with significant anatomical differences between a group of patients and a control group. Several image preprocessing steps are included in the method's pipeline and deformable image registration is the one playing a crucial role. Its central idea is to find local forces which will deform a floating image to make it more similar to a reference image. The involved non-linear transforms are either based on smooth basis functions [3], [4], [5] or they are physically interpreted, e.g. by mechanics of continuum [6], [7], [8]. The former group of methods produces smooth low-dimensional deformations which are able to suppress only gross anatomical inter-subject differences, whereas the goal of the latter methods is to achieve a perfect match. In [9], [10], images in VBM are put into a stereotaxic space by an affine transform and then they are warped to the reference image by low-dimensional parametric deformations based on lowest-frequency components of discrete cosine transform. The coefficients are searched in an optimization algorithm which minimizes the residual squared difference between the images and simultaneously maximizes the smoothness of the deformations. Only one scaling parameter is incorporated to count for differences in intensities of the images, what makes it suitable for monomodal images only.

In this paper, we propose a deformable registration algorithm proper for multimodal images. Below, we first explain the methods used in our algorithm and then we present experimental results obtained from its evaluation.

## Methods

Our deformable registration is performed by a multilevel block-matching technique, see fig 1. A floating image  $N$  is deformed to match a reference image  $M$  in an iterative process. A resulting displacement  $\mathbf{u}$  is made up from local translations of the blocks of the floating image  $N$  by radial basis function (RBF) interpolation. The translations representing warping forces  $\mathbf{f}$  are found by maximizing symmetric regional similarity measures. The floating image  $N$  is assumed to be brought into the same coordinate space as the reference image  $M$  by a previous linear registration step.

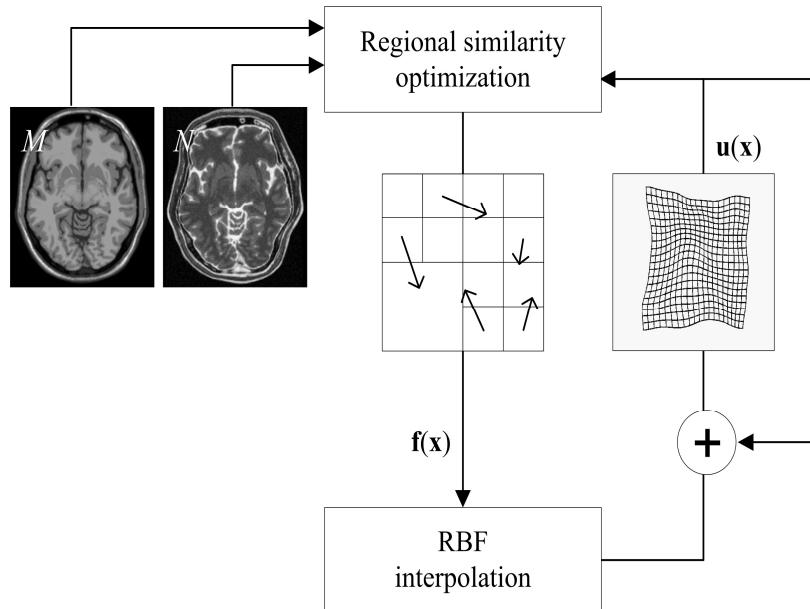


Fig. 1. Deformable registration scheme (see the text for details).

### Regional symmetric similarity measures

Various multimodal similarity measures are examined here. Regional similarity is computed by simply averaging point similarities over the region [11]:

$$S_w(\mathbf{w}) = \frac{1}{K_w} \sum_{\mathbf{x} \in W} S(\mathbf{x}) \quad (1)$$

where  $S_w$  denotes a similarity measure of a region  $W$  with the center point  $\mathbf{w}$  and  $K_w$  overlapping voxels  $\mathbf{x}$  in which point similarities  $S$  are computed. The point similarity measure  $S_{MI}$  derived from the well known global similarity measure mutual information is defined by [11]:

$$S_{MN}(x) = \log_2 \frac{p_{MN}(m(x), n(x))}{p_M(m(x)) \cdot p_N(n(x))}, \quad (2)$$

where  $p_{MN}$  denotes joint distribution of intensities and  $p_M$ ,  $p_N$  are marginal intensity distributions of the images  $M$  and  $N$  respectively. Another point similarity measure  $S_{UH}$  is proposed in [6]:

$$S_{UH}(x) = S_H(x) + S_{MN}(x) = \log_2 p_{MN}(m(x), n(x)) + S_{MN}(x), \quad (3)$$

where  $S_H$  is a point similarity measure derived from the global joint entropy of the images. All the defined point similarity measures depend on the joint intensity distribution, which is estimated from the joint histogram, which is not known until the images are perfectly matched. Thus, it is usually estimated from the images aligned by the previous registration step. In this way, the deformable registration is done also in [12], where a region similarity measure based on conditional probabilities is proposed. It is rewritten here as another point similarity measure:

$$S_{PC}(x) = p(n(x|m(x))), \quad (4)$$

which is defined by the probability of a correspondence between a given intensity  $m$  of the reference image  $M$  and any intensity  $n$  of the floating image  $N$ . The conditional probability densities are extracted by normalizing the values of each row of the joint histogram parallel to the axis with the intensities of the floating image  $N$ . Another similarity measure depending on probability rather than uncertainty is derived here from (2):

$$S_{PMI}(x) = \frac{p_{MN}(m(x), n(x))}{p_M(m(x)) \cdot p_N(n(x))}. \quad (5)$$

At each level of subdivision, translations of rectangular blocks of the floating image  $N$  are searched in an optimization algorithm, which maximizes a selected region similarity measure. Inspired by symmetric registration proposed in [13], the symmetric regional similarity measure is obtained here as a sum of two partial similarity measures. These are computed in the blocks of the floating image according to the reference image blocks as well as in the reverse direction.

To avoid getting trapped in local minima, a combination of extensive search and hillclimbing algorithms is used here. First, a space of all possible translations is searched with a relatively large step.  $P$  best points are then used as start points for the following hillclimbing. The minimum of  $P$  local minima obtained by the hillclimbing is then declared as the global minimum.

## Multilevel deformation

Once the local translations are found, the displacement  $\mathbf{u}$  is computed separately for each of  $D$  dimensions by interpolation with the use of RBF by:

$$u_k(\mathbf{x}) = \sum_{i=1}^B (\alpha_i \cdot R(\|\mathbf{x} - \mathbf{w}_i\|)), \quad k = 1 \dots D, \quad (6)$$

where  $u_k(\mathbf{x})$  is the displacement of a grid point  $\mathbf{x}$  in the  $k^{\text{th}}$  dimension,  $R$  is the radial basis function of the distance  $\|\mathbf{x} - \mathbf{w}_i\|$  between the grid point  $\mathbf{x}$  and the center of the  $i^{\text{th}}$  block  $\mathbf{w}_i$ . The coefficients  $\alpha_i$  are computed by putting the translations  $\mathbf{f}$  into (6) and solving the resulting linear system of  $B$  equations separately for each dimension  $k$ . The compactly supported Wendland's RBF, which was successfully used for landmark-based deformable registration in [4] is used here. Its mathematical properties hold for different spatial support, which is important for the multilevel strategy. For each level of subdivision, the block size is set to the half of the size at the previous level. The displacements are gradually incremented over all levels, refining the resulting deformation in the coarse-to-fine manner. The regions containing poor contour or surface information can be eliminated from the matching process and the algorithm can be accelerated in this way. The subdivision is performed only if at least one voxel in the current region has its normalized gradient image intensity bigger than a certain threshold, see fig. 2.

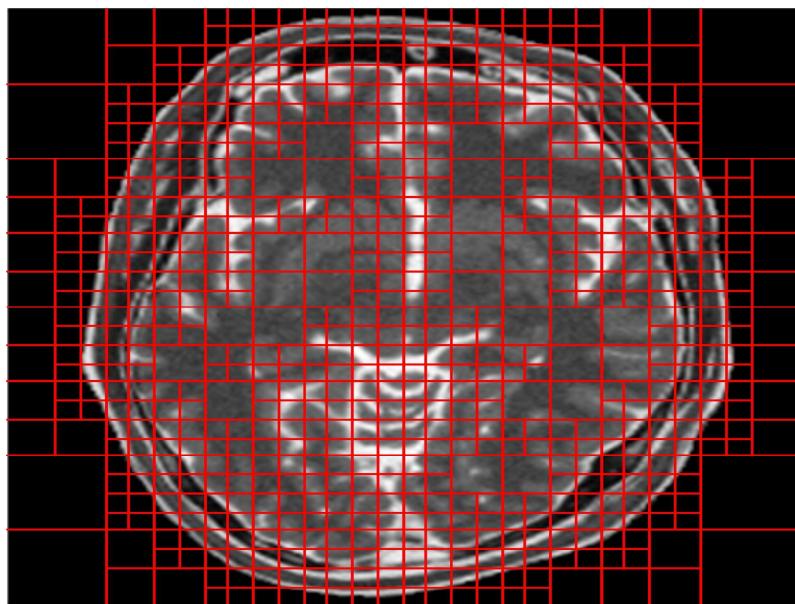


Fig. 2. Illustration of five-level adaptive subdivision.

## Tissue probability maps

The computation of the global joint histogram is not the only way how to estimate the joint intensity distribution of the images  $M$  and  $N$ . When the registration is done in a stereotaxic

Low-dimensional Multimodal Deformable Registration of MRI Brain Images in Stereotaxic Space

space, tissue probability maps are often available, representing a kind of prior information, which can be used. Here, another estimate of the joint intensity distribution is made with their use. It is further combined with the usual estimate obtained from the global joint histogram by a weighting parameter  $\lambda$ :

$$P_{\text{M&V}}(m, n) = \lambda \cdot P_{\text{M&V}}^{\text{prior}}(m, n) + (1 - \lambda) \cdot P_{\text{M&V}}^{\text{global}}. \quad (7)$$

The intensities representing main tissues in the brain images are emphasized in this way.

The joint intensity estimates are re-computed in each iteration of the registration algorithm. As the intensities of the floating image are not spread on a regular grid, simple unitary increments of individual histogram bins cannot be done. Instead of that, all histogram bins corresponding to the reference image voxels in the neighborhood of a displaced floating image voxel are increased by a value equal to the value of an interpolation kernel function. Cubic B-splines are used for that purpose in the generalized partial volume interpolation (GPVE) described in [14]. This approach is used here to compute the global joint histogram. In the computation of the joint distribution estimate based on the tissue probability maps, the interpolation has to be done among the voxels of the floating image. Thus, some ideas from the distance-weighted scattered data interpolation methods were adapted here and the locally bounded kernel function described in [15] is used in the partial volume interpolation scheme.

## Experiment and results

The performance of the proposed algorithm with various similarity measures was evaluated with the use of 2D image data obtained from Simulated Brain Database (SBD) [16]. The original size of the SBD transversal slices is  $181 \times 217$  pixels with the pixel size  $1 \times 1$  mm. For the evaluation purposes, the images were padded to the size of  $217 \times 217$  pixels. The square image size is convenient for the subdivision scheme. Synthetic deformations were composed from random translations at 10% randomly selected pixels in transversal slices. These force fields were smoothed by gaussian filtering with random standard deviation ( $\sigma = 10 \pm 5$  mm) to obtain final displacements which were then applied on 20 intensity images and corresponding segmented images. The average initial overlap error was 41.0%. The deformations were then recovered by the proposed deformable registration algorithm and the average decrease of the overlap error  $\Delta e$  was computed. T1 weighted, T2-weighted and PD (proton density) images with 3% noise and 20% intensity nonuniformity were used as floating images and the T1-weighted image with no noise and no intensity nonuniformity was used as the reference image. The results for various levels of decomposition are summarized in tab 1. The overlap error got smaller up to the 5<sup>th</sup> level, when the subimages had size of 7x7 pixels. Although the next level gave an increase in the global mutual information, the alignment measured by overlap errors and also by visual inspection was constant or worse. This level was considered as the maximum subdivision level for this algorithm.

Low-dimensional Multimodal Deformable Registration of MRI Brain Images in Stereotaxic Space

*Tab. 1. Quantitative validation of the proposed algorithm. The averages of the overlap error decrease  $\Delta e$  were computed for various similarity measures, various image pairs and various levels of decomposition.*

Images	$\Delta e [\%]$			
	$S_{PC}$	$S_{UH}$	$S_{MI}$	$S_{PMI}$
1 <sup>st</sup> level				
T1-T1	<b>5.2</b>	1.6	1.7	1.5
T1-T2	<b>5.3</b>	1.3	1.4	2.8
T1-PD	<b>4.4</b>	1.5	1.3	2.2
2 <sup>nd</sup> level				
T1-T1	<b>12.6</b>	8.3	8.4	8.5
T1-T2	<b>12.3</b>	4.7	4.8	7.4
T1-PD	<b>10.1</b>	5.6	6.0	5.8
3 <sup>rd</sup> level				
T1-T1	<b>22.6</b>	17.9	17.9	19.5
T1-T2	<b>21.2</b>	9.8	9.8	16.4
T1-PD	<b>17.1</b>	11.4	12.3	12.2
4 <sup>th</sup> level				
T1-T1	<b>27.2</b>	24.8	24.9	26.6
T1-T2	<b>25.3</b>	16.0	16.2	23.3
T1-PD	<b>20.2</b>	15.6	16.4	17.4
5 <sup>th</sup> level				
T1-T1	27.7	26.8	27.1	<b>28.5</b>
T1-T2	<b>25.1</b>	19.4	19.9	24.9
T1-PD	<b>20.0</b>	16.9	17.9	19.2

## Conclusion

An algorithm for low-dimensional atlas-based registration of MRI images was presented. Four various symmetric region similarity measures were studied in an experiment in which synthetic deformations were recovered and the performance of the algorithm was quantified by the decrease of overlap error in segmented images. The similarities were measured with the use of joint histogram and tissue probability maps from MRI brain atlases. The overlap error was lower with the similarity measures  $S_{PC}$  and  $S_{PMI}$  depending on probabilities than when the similarity measures  $S_{MI}$  and  $S_{UH}$  depending on uncertainty were used. In our implementation, partial volume interpolation scheme was used, so that it was unnecessary to compute the deformed floating image during the registration process. The proposed algorithm is suitable for the voxel based morphometry, as the precision of the registration can be controlled by the maximum level of decomposition. Thus, only gross inter-subject anatomical differences can be suppressed and the important variability for statistical parametric tests can

be preserved. In addition, the use of multimodal similarity measure allows to use an arbitrary available brain atlas without any need to transform intensities in the images to obtain monomodal data.

## Acknowledgement

This work has been partly supported by the grant projects 1ET2085120511 AC CR, 102/04/0472 and 305/04/1385 from GACR, and the Research Programme of Brno University of Technology MSM 0021630513.

## References

- [1] Friston, K. J., Ashburner,J.: Generative and Recognition Models for Neuroanatomy. *NeuroImage*, vol. 23, pp. 21–24, 2004.
- [2] Davatzikos, Ch.: Why Voxel-based Morphometric Analysis Should Be Used with Great Caution When Characterizing Group Differences. *NeuroImage*, vol. 23, pp. 17–20, 2004.
- [3] Pauchard, Y., Smith M. R., Mintchev M. P.: Modeling Susceptibility Difference Artifacts Produced by Metallic Implants in Magnetic Resonance Imaging with Point-Based Thin-Plate Spline Image Registration. In Proc. 26th Conf. IEEE EMBS, 2004, pp. 1766–1769.
- [4] Fornefett, M., Rohr, K., Stiehl, H. S.: Radial Basis Functions with Compact Support for Elastic Registration of Medical Images. *Image Vision Comput.*, vol. 19, pp. 87–96, 2001.
- [5] Rohlffing, T., Maurer, C. R., Bluemke D. A., Michael J. A.: Volume-Preserving Nonrigid Registration of MR Breast Images Using Free-Form Deformation With an Incompresibility Constraint. *IEEE Trans. on Medical Imaging*, vol. 12, pp. 730–741, 2003.
- [6] Rogelj, P., Kovačič, S., Gee, J. C.: Point Similarity Measures for Non-rigid Registration of Multi-modal Data. *Computer Vision and Image Understanding*, vol. 92, pp. 112–140, 2003.
- [7] Tang, S., Jiang, T.: Fast Nonrigid Medical Image Registration by Fluid Model. In: Proc. 6th Asian Conference on Computer Vision, 2004.
- [8] Christensen, G. E., He, J.: Consistent Nonlinear Elastic Image Registration. In: IEEE Proc. MMBIA, 2001, pp. 37–43.
- [9] Ashburner, J., Friston, K. J.: Voxel-Based Morphometry – The Methods. *NeuroImage*, vol. 11, pp. 805–821, 2000.
- [10] Ashburner, J., Friston, K. J.: Nonlinear Spatial Normalization Using Basis Functions. *Human Brain Mapping*, vol. 7, pp. 254–266, 1999.
- [11] Rogelj, P., Kovačič, S.: Point Similarity Measure Based on Mutual Information. In *Biomedical Image Registration: revised papers*, 2003, pp.112–121.
- [12] Maintz, J. B. A., Meijering, E. H. W., Viergever M. A.: General Multimodal Elastic Registration based on Mutual Information. In *Medical Imaging: Image Processing* (Proc. of SPIE, vol. 3338), 1998, pp. 144–154.

*Low-dimensional Multimodal Deformable Registration of MRI Brain Images in Stereotaxic Space*

- [13] Rogelj, P., Kovačič, S.: Symmetric Image Registration. In Medical Imaging: Image Processing (Proc. of SPIE, vol.5032), 2003, pp. 334–343.
- [14] Chen, H., Varshney, P. K.: Mutual Information-Based CT-MR Brain Image Registration Using Generalized Partial Volume Point Histogram Estimation. IEEE Trans. on Medical Imaging, vol. 22, pp. 1111–1119, 2003.
- [15] Franke, R. Nielson, G.: Smooth Interpolation of Large Sets of Scattered Data. International Journal for Numerical Methods in Engineering, vol. 15, pp. 1691–1704, 1980.
- [16] Collins, D.L., Zijdenbos, A.P., Kollokian, V., Sled, J.G., Kabani, N.J., Holmes, C. J., Evans A. C.: Design and Construction of a Realistic Digital Brain Phantom. IEEE Trans. on Medical Imaging, vol. 17, pp. 463–468, 1998.

# Málorozměrná multimodální pružná registrace obrazů mozku z MRI ve stereotaktickém prostoru

Daniel Schwarz<sup>1</sup>, Ivo Provažník<sup>2</sup>

1. Institute of Biostatistics and Analyses, Masaryk University, Czech Republic,

2. Department of Biomedical Engineering, Brno University of Technology, Czech Republic

Pružná registrace obrazů je klíčovou technikou ve výpočetní neuroanatomii. V tomto článku je navržena registrace metodou iterativního lícování podobrazů s několika vylepšením z oblasti mnohorozměrné multimodální registrace a registrace založené na význačných bodech. Je zde použita symetrická multimodální podobnostní metrika, díky které je možno registrovat obrazy různých subjektů na libovolný digitální atlas mozku. Výsledné deformace jsou získány pomocí interpolační techniky založené na radiálních bázových funkcích. Deformace potlačují v obrazech pouze hrubé tvarové rozdíly mezi různými subjekty a ponechávají v nich jemnou anatomickou variabilitu, která bývá předmětem následného zkoumání vybranými morfometrickými metodami. Součástí příspěvku jsou kvantitativní výsledky z experimentálního ověření navržené metody se čtyřmi různými podobnostními metrikami na obrazových datech ze simulátoru obrazů mozku.

**Klíčová slova:** zpracování obrazů, registrace obrazů, MRI, výpočetní neuroanatomie, radiální bázové funkce

## Úvod

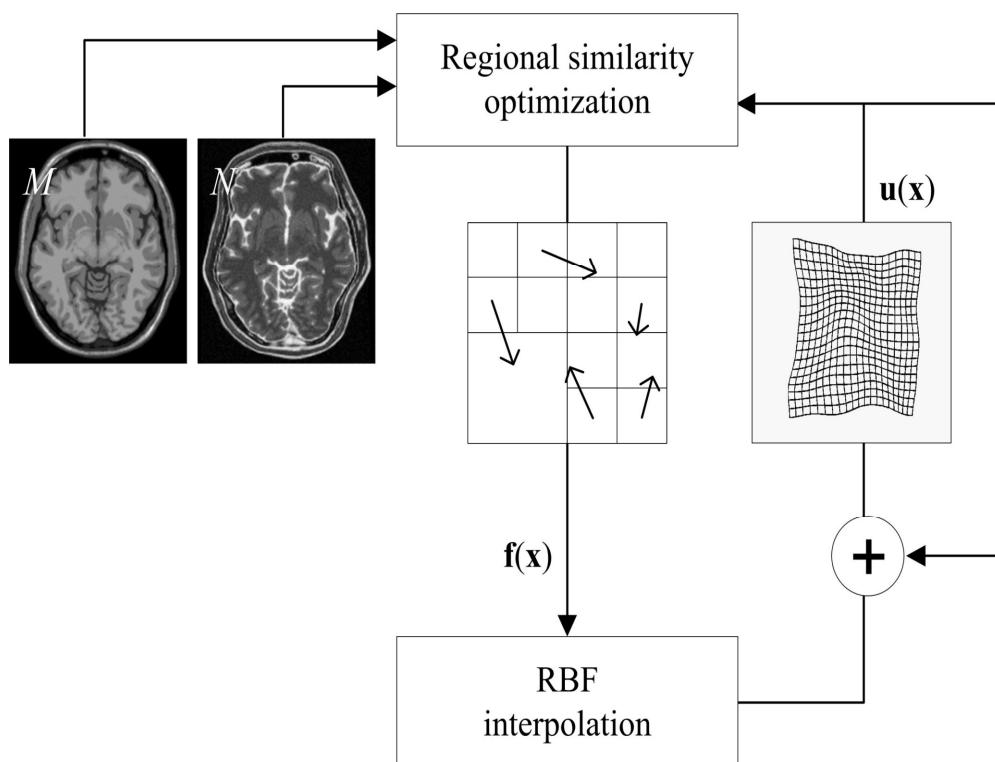
Jednou z intenzivně užívaných metod v oblasti výpočetní neuroanatomie je morfometrie založená na voxelech (voxel-based morphometry VBM), která je stále objektem aktivního výzkumu a diskuzí [1], [2]. Princip této metody spočívá ve zkoumání anatomických skenů z MRI voxel po voxelu s cílem automaticky vyznačit v mozku oblasti se signifikantními rozdíly mezi skupinou pacientů a kontrolní skupinou dobrovolníků. VBM je zřetězením několika algoritmů pro zpracování obrazů, přičemž samotná pružná registrace zde hraje klíčovou roli. Hlavní myšlenkou pružné registrace je nalézt lokální síly, které zdeformují plovoucí obraz tak, aby se více podobal obrazu referenčnímu. Použité nelineární transformace bývají založeny buď na hladkých bázových funkcích [3], [4], [5] nebo vychází z fyzikálních interpretací, např. z mechaniky těles [6], [7], [8]. Zatímco v prvním případě se získávají hladké málorozměrné deformace, kterými je možno potlačit pouze globální rozdíly ve tvaru mozku mezi jednotlivými subjekty, cílem druhé skupiny metod je dosáhnout v registrovaných obrazech perfektního slícování. V původní VBM [9], [10] jsou obrazy afinní transformací převedeny do stereotaktického prostoru a následně jsou zdeformovány na referenční obraz pomocí málorozměrné parametrické transformace založené na bázových funkcích užívaných v diskrétní kosinové transformaci. Koeficienty transformace jsou hledány optimalizačním algoritmem, který minimalizuje sumu čtverců rozdílů mezi intenzitami v obrazech a zároveň maximalizuje hladkost dosažených deformací. Rozdíly v intenzitách obrazů jsou postiženy

pouze jedním škálovacím parametrem, a tak je tato metoda vhodná pouze pro monomodální obrazy.

V tomto článku je navržen algoritmus pružné registrace vhodný pro multimodální obrazy. Jsou zde jednak vysvětleny použité metody a dále jsou prezentovány experimentální výsledky získané z pokusů při kvantitativním hodnocení tohoto algoritmu.

## Metody

Zde navržená pružná registrace se provádí pomocí víceúrovňového lícování podobrazů, viz obr. 1. Plovoucí obraz  $N$  je v iteračním procesu deformován tak, aby lícoval s referenčním obrazem  $M$ . Výsledné vychýlení  $\mathbf{u}$  je pro každý voxel interpolováno z lokálních translací v obrazu  $N$  s využitím radiálních bázových funkcí (RBF). Translace podobrazů reprezentující deformační síly  $\mathbf{f}$  jsou nalezeny v lokálních registracích, ve kterých je maximalizována symetrická regionová podobnostní metrika. Předpokládá se, že plovoucí obraz  $N$  byl v předešlém kroku transformován do souřadného systému obrazu  $M$  pomocí lineární registrace.



Obr. 1. Schéma pružné registrace (detaily jsou popsány v textu).

### Regionové symetrické podobnostní metriky

V navrženém algoritmu mohou být použity různé podobnostní metriky. Regionová podobnost se vypočítá zprůměrováním bodových podobností v podobrazu (dále regionu) [11]:

Málorozměrná multimodální pružná registrace obrazů mozku z MRI ve stereotaktickém prostoru

$$S_W(\mathbf{w}) = \frac{1}{K_W} \sum_{\mathbf{x} \in W} S(\mathbf{x}), \quad (1)$$

kde  $S_W$  označuje podobnostní metriku regionu  $W$  se středem v bodě  $\mathbf{w}$  a s  $K_W$  překrývajícími se voxely  $\mathbf{x}$ , ve kterých jsou vyhodnocovány bodové podobnosti  $S$ . Bodová podobnostní metrika  $S_{MI}$  odvozená z dobře známé globální podobnostní metriky zvané vzájemná informace je definována jako [11]:

$$S_{MI}(\mathbf{x}) = \log_2 \frac{p_{MN}(m(\mathbf{x}), n(\mathbf{x}))}{p_M(m(\mathbf{x})) \cdot p_N(n(\mathbf{x}))}, \quad (2)$$

kde  $p_{MN}$  označuje sdruženou hustotu pravděpodobnosti intenzit a  $p_M$ ,  $p_N$  jsou marginální hustoty pravděpodobnosti intenzit v obrazech  $M$  a  $N$ . Jiná bodová podobnostní metrika  $S_{UH}$  je navržena v [6]:

$$S_{UH}(\mathbf{x}) = S_H(\mathbf{x}) + S_{MI}(\mathbf{x}) = \log_2 p_{MN}(m(\mathbf{x}), n(\mathbf{x})) + S_{MI}(\mathbf{x}), \quad (3)$$

kde  $S_H$  je bodová podobnostní metrika odvozená z globální sdružené entropie obrazů. Všechny zde definované metriky závisejí na sdružené hustotě pravděpodobnosti, která bývá odhadována ze vzájemného histogramu obrazů, který ovšem není znám, dokud nejsou obrazy přesně slícovány. Sdružená hustota pravděpodobnosti proto bývá odhadována z obrazů slícových pomocí předchozí lineární registrace. Takto je tomu i ve [12], kde je navržena pro nelineární registraci regionová metrika založená na podmíněných pravděpodobnostech. Její definice je zde přepsána do podoby bodové podobnostní metriky:

$$S_{PC}(\mathbf{x}) = p(n(\mathbf{x}) \| m(\mathbf{x})), \quad (4)$$

která je tedy dána pravděpodobností závislosti mezi danou intenzitou  $m$  v referenčním obrazu  $M$  a intenzitou  $n$  plovoucího obrazu  $N$ . Hustoty podmíněných pravděpodobností jsou získány normalizací hodnot v každém řádku vzájemného histogramu rovnoběžném s osou intenzit plovoucího obrazu  $N$ . Poslední zkoumaná podobnostní metrika je zde odvozena z (2):

$$S_{PMI}(\mathbf{x}) = \frac{p_{MN}(m(\mathbf{x}), n(\mathbf{x}))}{p_M(m(\mathbf{x})) \cdot p_N(n(\mathbf{x}))}. \quad (5)$$

V každé úrovni dělení obrazů jsou hledána posunutí podobrazů v plovoucím obrazu  $N$  pomocí optimalizačního algoritmu, který maximalizuje vybranou regionovou podobnostní metriku. S využitím myšlenky symetrické registrace [13] je zde navržena symetrická regionová podobnostní metrika, která se počítá jako součet dvou metrik: pro přímou registraci podobrazu v  $N$  vzhledem k referenčnímu obrazu  $M$  a pro směr opačný.

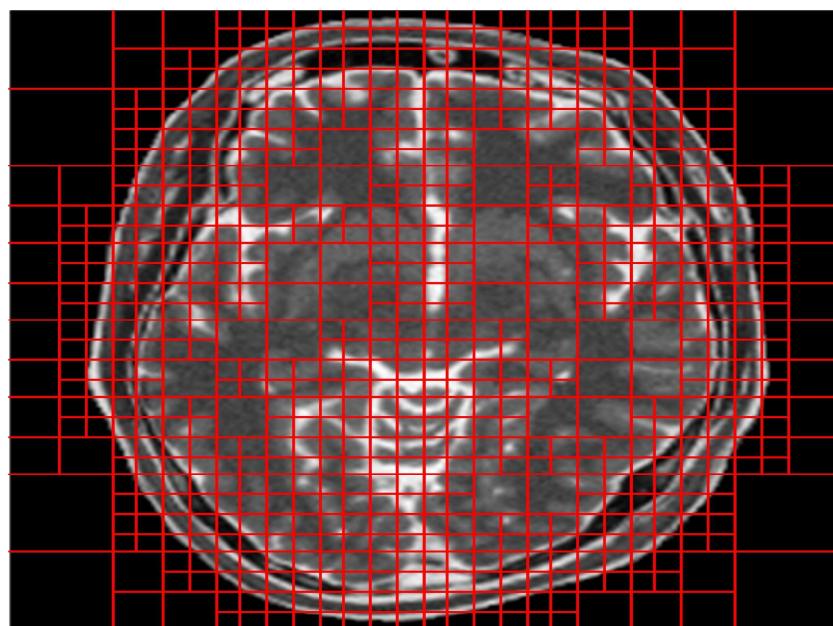
Optimalizační algoritmus je zde navržen s cílem vyhnout se lokálním optimům. Jedná se o kombinaci algoritmu rozsáhlého vyhledávání (z angl. extensive search) a algoritmu nejstrmějšího sestupu. V prvním kroku je prostor všech možných translací prohledán s relativně velkým krokem.  $P$  nejlepších bodů je pak určeno jako počáteční body pro algoritmus nejstrmějšího sestupu. Za globální maximum se vezme maximum z  $P$  lokálních maxim nalezených nestrmějším sestupem.

## Víceúrovňová deformace

Po nalezení lokálních translací všech podobrazů je spočítáno vychýlení  $\mathbf{u}$  interpolací s využitím RBF. Vychýlení je počítáno zvlášť pro každou z  $D$  souřadnicových os:

$$u_k(\mathbf{x}) = \sum_{i=1}^B (\alpha_i \cdot R(\|\mathbf{x} - \mathbf{w}_i\|)), \quad k = 1 \dots D, \quad (6)$$

kde  $u_k(\mathbf{x})$  je vychýlení bodu  $\mathbf{x}$  obrazové mřížky podél osy  $k$ ,  $R$  je radiální bázová funkce vzdálenosti  $\|\mathbf{x} - \mathbf{w}_i\|$  mezi bodem obrazové mřížky a středem podobrazu  $\mathbf{w}_i$ . Koeficienty  $\alpha_i$  se spočítají dosazením lokálních translací  $\mathbf{f}$  do (6) a vyřešením vzniklého systému  $B$  lineárních rovnic – opět pro každou souřadnicovou osu zvášť. Jako funkce  $R$  je zde použito Wendlandovy funkce s kompaktním nosičem, která byla úspěšně použita v [4] pro pružnou registraci založenou na význačných bodech. Její výhodné matematické vlastnosti umožňují její prostorovou dilataci a kompresi, což je důležité pro zavedení víceúrovňové strategie. Pro každou úroveň dělení podobrazů je nastavena velikost podobrazu na polovinu velikosti z předcházející úrovně. Vychýlení jsou sčítána přes všechny úrovně, čímž dochází k postupnému zpřesňování deformace. Algoritmus lze urychlit vynecháním z lícování těch oblastí, které neobsahují žádné kontury nebo povrchy. Další dělení podobrazu je prováděno jen v případě, když alespoň v jednom jeho voxelu je normalizovaný gradient obrazové intenzity větší než předdefinovaný práh, viz obr. 2.



Obr. 2: Pěti-úrovňové adaptivní dělení podobrazů.

## Tkáňové pravděpodobnostní mapy

Výpočet vzájemného histogramu obrazů není jediným způsobem, jak lze odhadnout sdruženou hustotu pravděpodobnosti obrazů  $M$  a  $N$ . Provádí-li se registrace ve stereotaktickém prostoru, často jsou k dispozici tkáňové pravděpodobnostní mapy reprezentující užitečnou apriorní informaci. S využitím této informace je zde počítán ještě jeden odhad sdružené hustoty pravděpodobnosti a tento je dále kombinován s obvyklým odhadem založeným na vzájemném histogramu obrazů váhovacím parametrem  $\lambda$ :

$$p_{MN}(m,n) = \lambda \cdot p_{MN}^{prior}(m,n) + (1-\lambda) \cdot p_{MN}^{global}. \quad (7)$$

Takto jsou ve výsledném odhadu  $p_{MN}$  zvýrazněny ty intenzitní páry, které reprezentují hlavní mozkové tkáně.

Během registrace jsou odhady sdružené hustoty pravděpodobnosti přepočítávány v každé iteraci. Nelze přitom použít klasické inkrementování v histogramových kontejnerech, poněvadž plovoucí obraz nemá vzhledem k dosud vypočteným vychýlením pravidelnou mřížku. Je proto použito zobecněné interpolace částečných objemů (GPVE) [14], kdy histogramové kontejnery související se všemi voxely referenčního obrazu sousedícími s vychýleným voxellem plovoucího obrazu jsou navýšeny o hodnotu interpolační jádrové funkce, kterou je v tomto případě kubický B-spline. Tohoto postupu je zde využito pro výpočet vzájemného histogramu obrazů. Při výpočtu založeném na tkáňových pravděpodobnostních mapách se interpolace provádí mezi voxely plovoucího obrazu, a tak byl tento postup adaptován do interpolace nepravidelně rozložených dat založené na váhování vzdálenostmi [15].

## Experimentální výsledky

Kvalita registrace navrženým algoritmem s využitím různých podobnostních metrik byla vyhodnocena s využitím 2D obrazových dat získaných z databáze simulovaných obrazů (Simulated Brain Database SBD) [16]. Velikost transverzálních řezů v SBD obrazech je  $181 \times 217$  pixelů s velikostí pixelu  $1 \times 1$  mm. Obrazy byly rozšířeny na velikost  $217 \times 217$  pixelů, poněvadž čtvercový rozměr je výhodnější pro víceúrovňové dělení. Byly vygenerovány umělé deformace, a to z náhodných translací umístěných do 10 % náhodně vybraných pixelů. Tato silová pole byla vyhlazena pomocí gaussovských filtrov s náhodně volenou směrodatnou odchylkou ( $\sigma=10\pm5$  mm). Výsledná vychýlení byla aplikována na 20 intenzitních obrazů a k nim korespondující segmentované obrazy. Průměrná počáteční chyba překryvu segmentovaných částí v původních nezkreslených obrazech a obrazech po deformaci byla 41,0%. Deformace pak byly potlačeny pomocí zde navržené pružné registrace a byl vypočten průměrný pokles chyby překryvu  $\Delta e$ . Pro plovoucí obrazy byly požity obrazy váhované T1, T2 a PD (hustotou protonů). Tyto obrazy byly zarušeny šumem a artefaktem označovaným jako nelinearity přenosu obrazové informace (intensity nonuniformity). Jako referenční obraz posloužil T1 váhovaný obraz bez simulovaných artefaktů. Výsledky pro různé maximálně úrovně dekompozice jsou v tab. 1. Chyba překryvu se zmenšovala až do 5. úrovně dělení, kdy byla velikost podobrazů  $7 \times 7$  pixelů. Při dalším dělení podobrazů sice vzrostla ještě globální vzájemná informace, ovšem slícování obrazů hodnocené chybou

Málorozměrná multimodální pružná registrace obrazů mozku z MRI ve stereotaktickém prostoru

překryvu a vizuální kontrolou bylo stejné nebo horší. Proto je zde 5. úroveň považována za maximální úroveň dělení podobrazů pro tento regisrační algoritmus.

*Tab. 1. Kvantitativní hodnocení navrženého algoritmu pomocí poklesů průměrné chyby překryvu  $\Delta e$ . Jsou zde prezentovány výsledky získané při použití různých podobnostních metrik, pro různé páry intenzitních obrazů a pro různé úrovně dělení podobrazů.*

Obrazy	$\Delta e [\%]$			
	$S_{PC}$	$S_{UH}$	$S_{MI}$	$S_{PMI}$
1. úroveň				
T1-T1	<b>5,2</b>	1,6	1,7	1,5
T1-T2	<b>5,3</b>	1,3	1,4	2,8
T1-PD	<b>4,4</b>	1,5	1,3	2,2
2. úroveň				
T1-T1	<b>12,6</b>	8,3	8,4	8,5
T1-T2	<b>12,3</b>	4,7	4,8	7,4
T1-PD	<b>10,1</b>	5,6	6,0	5,8
3. úroveň				
T1-T1	<b>22,6</b>	17,9	17,9	19,5
T1-T2	<b>21,2</b>	9,8	9,8	16,4
T1-PD	<b>17,1</b>	11,4	12,3	12,2
4. úroveň				
T1-T1	<b>27,2</b>	24,8	24,9	26,6
T1-T2	<b>25,3</b>	16,0	16,2	23,3
T1-PD	<b>20,2</b>	15,6	16,4	17,4
5. úroveň				
T1-T1	27,7	26,8	27,1	<b>28,5</b>
T1-T2	<b>25,1</b>	19,4	19,9	24,9
T1-PD	<b>20,0</b>	16,9	17,9	19,2

## Závěr

V tomto článku byl navržen algoritmus pro málorozměrnou registraci MRI obrazů vhodnou pro pružnou registraci obrazů na digitální atlas mozku. V experimentu s umělými deformacemi byla studována kvalita registrace při použití čtyř různých podobnostních metrik. Kvalita registrace byla kvantitativně hodnocena pomocí poklesu chyby překryvu v segmentovaných obrazech. Podobnosti podobrazů byla měřeny s využitím vzájemného histogramu obrazů a dále také s využitím tkáňových pravděpodobnostních map. Chyba překryvu byla nejmenší při použití metrik  $S_{PC}$  a  $S_{PMI}$ . Při implementaci algoritmu bylo použito zobecněné interpolace částečných objemů, takže nebylo nutné během regisračního procesu přepočítávat deformovaný plovoucí obraz. Navržený regisrační algoritmus je vhodný pro

automatickou morfometrii v obrazech označovanou jako voxel-based morphometry, protože přesnost registrace lze zde řídit nastavením maximální úrovně dělení. Takto mohou být potlačeny v obrazech pouze globální tvarové rozdíly, zatímco jemná anatomická variabilita, která je důležitá pro statistické parametrické testy, může být uchována. Navíc díky použití multimodálních podobnostních metrik je možné tímto algoritmem registrovat obrazy na libovolný atlas mozku bez nutnosti další transformace intenzit v registrovaných obrazech.

## Poděkování

Tato práce byla podpořena granty 1ET2085120511 AV ČR, 102/04/0472 a 305/04/1385 GAČR a dále výzkumným záměrem Vysokého učení technického v Brně MSM 0021630513.

## Literatura

- [1] Friston, K. J., Ashburner, J.: Generative and Recognition Models for Neuroanatomy. *NeuroImage*, vol. 23, pp. 21–24, 2004.
- [2] Davatzikos, Ch.: Why Voxel-based Morphometric Analysis Should Be Used with Great Caution When Characterizing Group Differences. *NeuroImage*, vol. 23, pp. 17–20, 2004.
- [3] Pauchard, Y., Smith M. R., Mintchev M. P.: Modeling Susceptibility Difference Artifacts Produced by Metallic Implants in Magnetic Resonance Imaging with Point-Based Thin-Plate Spline Image Registration. In Proc. 26th Conf. IEEE EMBS, 2004, pp. 1766–1769.
- [4] Fornefett, M., Rohr, K., Stiehl, H. S.: Radial Basis Functions with Compact Support for Elastic Registration of Medical Images. *Image Vision Comput.*, vol. 19, pp. 87–96, 2001.
- [5] Rohlfing, T., Maurer, C. R., Bluemke D. A., Michael J. A.: Volume-Preserving Nonrigid Registration of MR Breast Images Using Free-Form Deformation With an Incompressibility Constraint. *IEEE Trans. on Medical Imaging*, vol. 12, pp. 730–741, 2003.
- [6] Rogelj, P., Kovačič, S., Gee, J. C.: Point Similarity Measures for Non-rigid Registration of Multi-modal Data. *Computer Vision and Image Understanding*, vol. 92, pp. 112–140, 2003.
- [7] Tang, S., Jiang, T.: Fast Nonrigid Medical Image Registration by Fluid Model. In: Proc. 6th Asian Conference on Computer Vision, 2004.
- [8] Christensen, G. E., He, J.: Consistent Nonlinear Elastic Image Registration. In: IEEE Proc. MMBIA, 2001, pp. 37–43.
- [9] Ashburner, J., Friston, K. J.: Voxel-Based Morphometry – The Methods. *NeuroImage*, vol. 11, pp. 805–821, 2000.
- [10] Ashburner, J., Friston, K. J.: Nonlinear Spatial Normalization Using Basis Functions. *Human Brain Mapping*, vol. 7, pp. 254–266, 1999.
- [11] Rogelj, P., Kovačič, S.: Point Similarity Measure Based on Mutual Information. In: Biomedical Image Registration: revised papers, 2003, pp. 112–121.
- [12] Maintz, J. B. A., Meijering, E. H. W., Viergever M. A.: General Multimodal Elastic Registration based on Mutual Information. In: Medical Imaging: Image Processing (Proc. of SPIE, vol. 3338), 1998, pp. 144–154.

Málorozměrná multimodální pružná registrace obrazů mozku z MRI ve stereotaktickém prostoru

- [13] Rogelj, P., Kovačič, S.: Symmetric Image Registration. In Medical Imaging: Image Processing (Proc. of SPIE, vol.5032), 2003, pp. 334–343.
- [14] Chen, H., Varshney, P. K.: Mutual Information-Based CT-MR Brain Image Registration Using Generalized Partial Volume Point Histogram Estimation. IEEE Trans. on Medical Imaging, vol. 22, pp. 1111–1119, 2003.
- [15] Franke, R. Nielson, G.: Smooth Interpolation of Large Sets of Scattered data. International Journal for Numerical Methods in Engineering, vol. 15, pp. 1691–1704, 1980.
- [16] Collins, D.L., Zijdenbos, A.P., Kollokian, V., Sled, J.G., Kabani, N.J., Holmes, C. J., Evans A. C.: Design and Construction of a Realistic Digital Brain Phantom. IEEE Trans. on Medical Imaging, vol. 17, pp. 463–468, 1998.

# A Prediction of Blood Flow through a Bypass Graft Using Statistical Methods

Jana Vrbková<sup>1</sup>, Vilém Bruk<sup>2</sup>

1. Department of Mathematical Analysis and Applications of Mathematics, Faculty of Natural Sciences, Palacky University Olomouc, Czech Republic,  
2. Clinic of Cardiosurgery, Faculty Hospital Olomouc, Czech Republic

Myocardial revascularization belongs among the most frequent cardiosurgery operations. Perioperative and longterm survival depend on the patency of the graft used and the anastomotic quality. Haemodynamical characteristics measured during a coronary artery bypass graft (CABG) surgery help verify anastomotic quality and also affect longterm graft patency. During CABG surgery (on a heart bypass machine with extracorporeal circulation), a surgeon measures blood flow through the bypass at the time the cross clamp is applied to the ascending aorta (blood is not flowing through coronary vessels, rather through the bypass) and later at the same place after removal of the cross clamp. The aim of this article is to find a statistical model for prediction of blood flow through the bypass after removal of the cross clamp based on the blood flow value when the cross clamp is placed on the aorta. When this prediction is good, we will be able to decrease a number of measurements with keeping whole information about an object.

**Keywords:** myocardial revascularization, prediction of blood flow and blood pressure, multiple regression, nonlinear regression model, linearization, linear regression model with constraints, outlier, leverage

## Introduction

Coronary revascularization belongs among the most frequent cardiosurgery operations in the world as well as in our republic (in the year 2002, 10797 cardiosurgery operations were performed of which 7051 were CABG surgeries) [1]. Perioperative and longterm survival depend on the patency of the graft used and the anastomotic quality. Many authors were concerned with finding characteristics which determine graft patency and anastomotic quality from both a short-term and a long-term point of view. As it turns out, haemodynamical characteristics help to verify anastomotic quality as well as longterm graft patency [2], [3], [4], [5], [6]. Louagie indicates resistance as a dominant characteristic for longterm graft patency. Resistance can be calculated as blood pressure divided by blood flow through a vessel [7]. Although this formula is only a simplification of the reality (blood flow is not a steady flow and blood is not a Newtonian fluid) [8], it is reasonable to think about blood pressure and blood flow as characteristics which determine resistance. Hata [9] confirmed on a set of patients with low free blood flow through an arterial bypass that a left mammary artery (LIMA) is able to adapt its diameter in relation to storage demands of a target coronary bed. It is also known that a flow through a bypass is affected by the competitive flow of a native coronary bed. Furthermore, flow through a bypass is speculated to affect changes of graft patency in relation to a competitive flow and the sensitivity of different kind of bypass

grafts (LIMA, RIMA, SVG) [10], [11], [12]. During myocardial revascularization (CABG on a heart bypass machine), various measurements of blood flow through the bypass are taken: free-flow (the bypass is not anastomosed on a target coronary bed yet), before removal of the clamp from the aorta (a coronary bed is stored only by bypass – there is no competitive flow) and after removal of the cross clamp from the aorta (with a competitive flow). Sometimes, a surgeon takes two measurements before removal of the clamp from the aorta: one before removal of clamps from other bypasses and one after removal of these clamps (blood flow through a bypass also depends on collaterals). A flowmeter (produced by the Norwegian company Medi-Stim) is used for the measurement of blood flow during aortocoronary bypass surgery at the Faculty Hospital in Olomouc. This machine works by the so called Transit-Time method. The principle is based on the fact that the time required for the ultrasound to pass through blood is slightly longer when it is passing upstream than downstream [13]. A surgeon can observe perioperatively with the aid of this machine the blood flow curve, the blood flow value, the mean flow, the Pulsatility Index [PI=(max. flow – min. flow)/mean flow] and other characteristics. It is possible to connect to this machine two flow probes, two pressure inputs and two auxiliary devices, for example ECG input. The flowmeter can calculate many characteristics from recorded data, for example the Fast Fourier Analysis for a saved curve (of blood flow, pressure etc.). In a small group of patients (35), mean arterial pressure and blood flow through a bypass (a left mammary artery grafted to the left anterior descending artery – LIMA-LAD) were measured during CABG surgery at the time when the cross clamp is applied to the aorta (blood is flowing to a coronary bed only through a bypass) and also at the time when the clamp is removed (a competitive flow). The aim of this article is to find a statistical model for the prediction of blood flow through a bypass after removal of the cross clamp from the aorta based on the blood flow value with the clamp still in the place on the aorta. When this prediction is good, it would be able to decrease a number of measurements with keeping whole information about an object.

## 1. Input data and model

### 1.1 Input data

Data representing blood flow through a bypass (LIMA-LAD) are used for the analysis. The measurements were taken two times: with the cross clamp on the aorta (LIMA-LAD I) and after the clamp was removed (LIMA-LAD III). At the same time, mean arterial pressure was recorded in mmHg and the flow in ml/min. Data were obtained from 35 patients during LIMA or BIMA (the measurements were taken on the left mammary artery) CABG surgeries.

Data obtained by measurement at the Clinic of Cardiosurgery, Faculty Hospital in Olomouc are organized in a vector of input data  $\zeta = (x_1, y_1, \dots, x_n, y_n, \xi_1, \eta_1, \dots, \xi_n, \eta_n)^T$  (' signs vector transposition), where

$x_1, x_2, \dots, x_n$	is blood flow through a bypass LIMA-LAD I,
$y_1, y_2, \dots, y_n$	is mean arterial blood pressure LIMA-LAD I,
$\xi_1, \xi_2, \dots, \xi_n$	is blood flow through a bypass LIMA-LAD III,
$\eta_1, \eta_2, \dots, \eta_n$	is mean arterial blood pressure LIMA-LAD III.

Table 1. Input data.

i	x <sub>i</sub>	y <sub>i</sub>	ξ <sub>i</sub>	η <sub>i</sub>
1	80	67	67	65
2	101	47	40	47
3	94	68	66	66
4	53	63	56	63
5	46	66	41	66
6	30	63	22	67
7	34	65	11	64
8	39	60	42	60
9	38	77	58	65
10	55	60	36	55
11	55	63	35	67
12	31	60	24	57
13	85	67	38	65
14	33	80	15	77
15	36	70	14	70
16	74	75	55	65
17	105	72	98	78

18	26	50	15	54
19	51	62	38	74
20	29	77	27	77
21	84	60	79	67
22	40	70	49	69
23	75	65	72	60
24	24	60	12	56
25	44	77	24	75
26	56	62	29	54
27	38	60	8	52
28	10	67	2	63
29	60	65	46	60
30	30	54	13	47
31	99	59	96	66
32	44	76	33	74
33	36	65	21	57
34	31	57	22	60
35	43	76	35	77

## 1.2 Statistical model

The model of the relationship between blood flow and pressure at the time when the clamp is applied to the aorta (LIMA-LAD I) and at the time when the clamp is removed (LIMA-LAD III) can be assumed in a few forms. The simplest possible form is the classical linear regression model.

$$\begin{pmatrix} v_{i,1} \\ v_{i,2} \end{pmatrix} = \begin{pmatrix} \square_1 \\ \square_2 \end{pmatrix} + \begin{pmatrix} \square_{11} & \square_{12} \\ \square_{21} & \square_{22} \end{pmatrix} \begin{pmatrix} x_i \\ y_i \end{pmatrix}, \quad (1)$$

where  $x_i$  and  $y_i$  are values of the patient's blood flow and blood pressure at the time when the cross clamp is applied to the aorta and  $v_{i,1}$  and  $v_{i,2}$  are predicted values of these parameters after removal of the cross clamp. There are random errors only on the left side in  $v_{i,1}$  and  $v_{i,2}$  values in this model. Parameters  $\alpha_1$ ,  $\alpha_2$ ,  $\beta_{11}$ ,  $\beta_{12}$ ,  $\beta_{21}$ ,  $\beta_{22}$  are estimated by the least square method.

However, the presumption of errorless values  $x_i$  and  $y_i$  don't correspond to the reality. The "Crystal ball" model, in which we assumed values  $x_i$  and  $y_i$  as realizations of random variables, is better.

$$E(\zeta) = [\zeta_{1,1}, \zeta_{1,2}, \dots, \zeta_{n,1}, \zeta_{n,2}, \zeta_{1,1}, \zeta_{1,2}, \dots, \zeta_{n,1}, \zeta_{n,2}], \quad Var(\zeta) = \begin{pmatrix} 2n & \mathbf{0} \\ \mathbf{0}^T & I \\ \end{pmatrix}_{2n, 2n}, \quad (2)$$

where  $E(\zeta)$  satisfies the conditions:

$$\begin{pmatrix} E(\xi_i) \\ E(\eta_i) \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} \beta_3 & \beta_4 \\ \beta_5 & \beta_6 \end{pmatrix} \begin{pmatrix} E(x_i) \\ E(y_i) \end{pmatrix}, \quad i=1,2,\dots,n. \quad (3)$$

Note:  $E(\zeta)$  and  $Var(\zeta)$  denote the mean value and the variance of the random vector  $\zeta$ .

The variance matrix  $Var(\zeta)$  is assumed in the simplest possible form and it is used for less complexity of following calculations. It will be necessary to assume more complex form of this matrix in further research but at the same time it will lead to a problem of estimation of variance components.

## 2. Prediction of blood flow and pressure

### 2.1 The classical regression model

Estimations of parameters  $\alpha_1, \alpha_2, \beta_{11}, \beta_{12}, \beta_{21}, \beta_{22}$ , which we obtain for our data set with usage of function lm(y~x)(available in software R), are shown in the table 2.

Table 2. Estimated parameters in the classical regression model.

$\hat{\alpha}$		$\hat{\beta}$	[,1]	[,2]
[1,]	-37.54	[1,]	0.82	0.044
[2,]	5.42	[2,]	0.52	0.86

In the figure 1 we can see the measured (black) and estimated (bold red) values of blood flow through a bypass and blood pressure after removal of the cross clamp from the aorta. To compare the efficiency of the classical regression model and the “Crystal ball” model we show residuals in the following table ( $\xi_i$  and  $\eta_i$  are measured values of blood flow and blood pressure after removal of the clamp from the aorta):

$$\Delta v_{i,1} = \xi_i - v_{i,1}, \quad \Delta v_{i,2} = \eta_i - v_{i,2}, \quad i=1,2,\dots,n. \quad (4)$$

A Prediction of Blood Flow through a Bypass Graft Using Statistical Methods

Table 3. The residuals in the classical regression model.

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
1	4.76	-1.69
2	-29.08	-3.36
3	-8.18	-2.16
4	17.85	0.94
5	7.01	1.66
6	2.62	5.95
7	-12.68	1.05
8	16.82	1.14
9	24.88	-8.47
10	-2.24	-4.56
11	-4.78	4.85

12	5.34	-1.51
13	-28.32	-1.91
14	-15.59	1.16
15	-13.89	2.65
16	-6.47	-8.32
17	12.78	5.91
18	5.57	4.33
19	1.99	12.89
20	1.22	3.92
21	17.10	6.17
22	17.85	1.47
23	14.87	-4.74

24	-0.94	-2.20
25	-14.02	1.26
26	-11.09	-7.33
27	-16.37	-6.82
28	-3.12	-0.63
29	1.11	-4.09
30	-1.75	-6.29
31	22.37	5.38
32	-4.50	1.13
33	-4.31	-6.04
34	4.89	4.08
35	-1.69	4.17

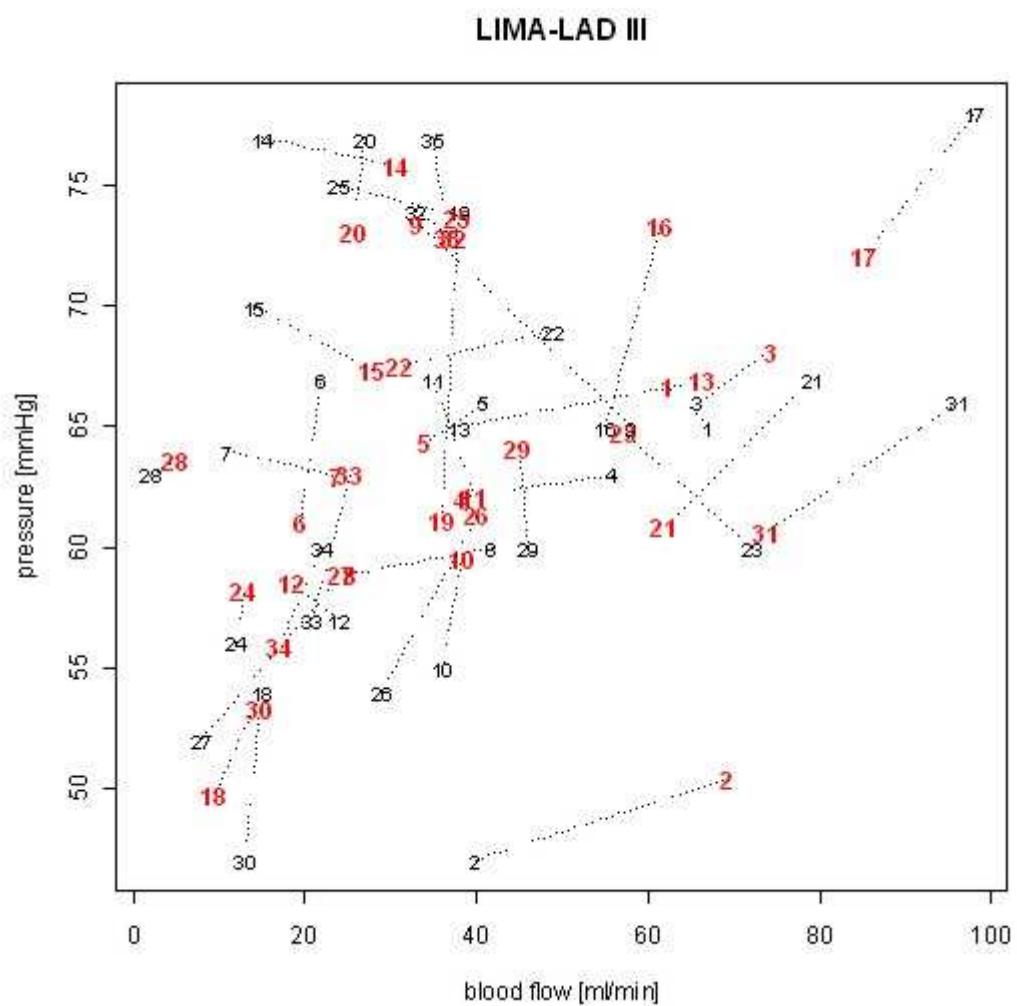


Fig. 1. Estimated and measured values of blood pressure and blood flow through a bypass after removal of the cross clamp from the aorta in the classical regression model.

## 2.2 The “Crystal ball” model

We establish predicted values for patient number i as follows:

$$\begin{pmatrix} \hat{v}_{i,1} \\ \hat{v}_{i,2} \end{pmatrix} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} + \begin{pmatrix} \hat{\beta}_3 & \hat{\beta}_4 \\ \hat{\beta}_5 & \hat{\beta}_6 \end{pmatrix} \begin{pmatrix} x_i \\ y_i \end{pmatrix}, \quad (5)$$

where  $x_i$  and  $y_i$  are measured values of blood flow and blood pressure at the time when the cross clamp is in the place on the aorta and  $\hat{v}_{i,1}$  and  $\hat{v}_{i,2}$  are predicted values of these parameters after removal of the cross clamp.

An algorithm of the calculation of estimations of parameters  $\beta_i$ ,  $i=1,2,\dots,6$ , which is iterative, is presented in the appendix A.1. Now we are testing the efficiency of the suggested algorithm on a data set which was used for the estimate of parameters  $\beta_i$ . We compare estimated values with directly measured values. The points for the initial iteration lie on a border of the data field ( $i=14, 17, 18$ ). The numbers mean order numbers of observations (patients).

Values of parameters  $\beta_i$  obtained in six iterative steps for the points of initial iteration mentioned above and our data set are shown in the table 4.

*Table 4. Estimated parameters in the “Crystal ball” model.*

$\hat{\beta}_1$	-99.10
$\hat{\beta}_2$	-11.96
$\hat{\beta}_3$	0.96
$\hat{\beta}_4$	1.34
$\hat{\beta}_5$	0.06
$\hat{\beta}_6$	1.12

For our data set we obtain following variance matrix and estimation of  $\sigma^2$  (the algorithm is presented in the appendix A.1):

*A Prediction of Blood Flow through a Bypass Graft Using Statistical Methods*

$$\hat{\sigma}^2 = 46.52034,$$

$$Var \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \\ \hat{\beta}_4 \\ \hat{\beta}_5 \\ \hat{\beta}_6 \end{pmatrix} = \begin{pmatrix} 539.914 & 225.030 & -0.840 & -7.528 & -0.350 & -3.137 \\ 225.030 & 326.432 & -0.350 & -3.137 & -0.508 & -4.551 \\ -0.840 & -0.350 & 0.008 & 0.005 & 0.003 & 0.002 \\ -7.528 & -3.137 & 0.005 & 0.110 & 0.002 & 0.046 \\ -0.350 & -0.508 & 0.003 & 0.002 & 0.005 & 0.003 \\ -3.137 & -4.551 & 0.002 & 0.046 & 0.003 & 0.066 \end{pmatrix}.$$

For numerical expression of accuracy of estimates obtained in our model, we show values of the following residuals in the table 5:

$$\Delta v_{i,1} = \xi_i - \hat{v}_{i,1}, \quad \Delta v_{i,2} = \eta_i - \hat{v}_{i,2}, \quad i=1,2,\dots,n. \quad (6)$$

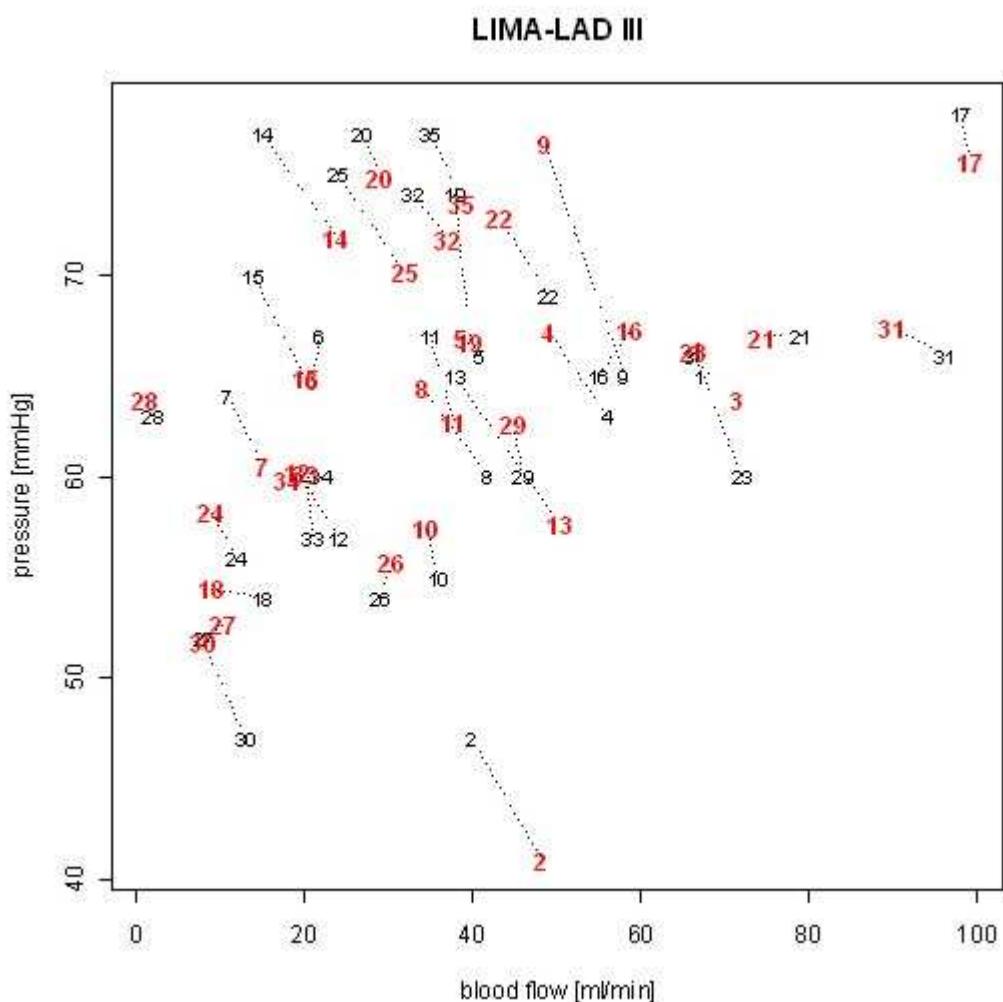
*Table 5. The residuals in the “Crystal ball” model.*

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
[1]	0.32	-1.33
[2]	-8.18	6.01
[3]	-5.39	2.25
[4]	6.95	-4.19
[5]	2.33	-0.94
[6]	1.14	2.23
[7]	-4.04	3.42
[8]	7.97	-4.39
[9]	9.40	-11.58
[10]	1.45	-2.46
[11]	-2.67	4.26

[12]	4.77	-3.28
[13]	-12.31	7.37
[14]	-8.76	5.08
[15]	-6.23	5.09
[16]	-3.81	-2.27
[17]	-1.16	2.36
[18]	6.04	-0.41
[19]	-1.75	7.33
[20]	-2.22	2.10
[21]	4.42	0.06
[22]	5.79	-3.86
[23]	5.66	-6.17

[24]	2.96	-2.28
[25]	-8.04	4.90
[26]	-1.69	-1.72
[27]	-2.38	-0.68
[28]	0.76	-0.76
[29]	1.11	-2.61
[30]	4.86	-4.78
[31]	6.06	-1.41
[32]	-4.19	2.26
[33]	0.84	-3.14
[34]	3.89	0.16
[35]	-3.87	3.39

In the figure 2, there are measured (black) and estimated (bold red) values of blood flow and blood pressure after removal of the cross clamp from the aorta.



*Fig. 2. Predicted and measured values of blood pressure and blood flow through a bypass after removal of the cross clamp in the “Crystal ball” model.*

Already from the picture and also from the values of residuals we can see that the “Crystal ball” model, which contains more parameters and this fact leads to better approximation of the measured values by our estimates, describes reality better than the classical linear regression model.

### 2.3 Location of outlier points in the “Crystal ball” model

In this section, we try to locate outlier observations in the “Crystal ball” model. Thereafter, we will exclude the outliers from the data set used in calculation of parameters and study changes in the model.

We find the indexes for which the expressions (35) are greater than or equal to 1.96 (or approximately 2). The points of these indexes are our suspect outliers.

For our data, set we have found as outliers the points with indexes 2, 9, 13, 14, 25, as you can see in the table 6.

A Prediction of Blood Flow through a Bypass Graft Using Statistical Methods

Table 6. Suspect outliers in the “Crystal ball” model.

i	$\{v_1\}_{\text{obs}} / \sqrt{\hat{Var}(v_1)_{\text{obs}}}$	$\{v_2\}_{\text{obs}} / \sqrt{\hat{Var}(v_2)_{\text{obs}}}$	i	$\{v_1\}_{\text{obs}} / \sqrt{\hat{Var}(v_1)_{\text{obs}}}$	$\{v_2\}_{\text{obs}} / \sqrt{\hat{Var}(v_2)_{\text{obs}}}$
1	-0.0620	0.0795	18	<b>-1.4524</b>	-0.0809
1	0.1952	-0.2574	19	0.3311	-0.4250
<b>2</b>	<b>2.4664</b>	<b>-2.4933</b>	19	<b>-1.0514</b>	1.3836
2	0.9999	1.4239	20	0.5530	-0.5654
3	1.3696	-1.3573	20	0.1288	0.4152
3	0.9032	0.4407	21	<b>-1.1524</b>	1.1128
4	-1.6880	1.6916	21	<b>-1.1412</b>	0.0126
4	-0.8654	-0.7923	22	<b>-1.4256</b>	1.4346
5	-0.5717	0.5661	22	-0.6574	-0.7431
5	-0.3838	-0.1771	23	<b>-1.3583</b>	1.4030
6	-0.3219	0.2789	23	-0.1483	-1.1888
6	-0.7385	0.4257	24	-0.7258	0.7348
7	0.9763	-0.9934	24	-0.2723	-0.4410
7	0.3050	0.6549	<b>25</b>	<b>1.9768</b>	<b>-1.9814</b>
8	-1.9367	1.9347	25	1.0032	0.9382
8	-1.0731	-0.8284	26	0.4584	-0.4183
<b>9</b>	<b>-2.2851</b>	<b>2.3823</b>	26	0.7779	-0.3309
<b>9</b>	0.0311	-2.2822	27	0.6324	-0.6015
10	-0.3321	0.3563	27	0.7532	-0.1346
10	0.1412	-0.4691	28	-0.1868	0.1920
11	0.6068	-0.6475	28	-0.0345	-0.1493
11	-0.2048	0.8037	29	-0.2397	0.2690
12	-1.1603	1.1688	29	0.2549	-0.4925
12	-0.5146	-0.6253	30	<b>-1.2062</b>	1.2371
<b>13</b>	<b>3.0490</b>	<b>-3.0550</b>	30	-0.2399	-0.9468
13	1.5730	1.4212	31	<b>-1.6281</b>	1.5954
<b>14</b>	<b>2.2011</b>	<b>-2.2022</b>	31	-1.3074	-0.2888
14	1.1691	0.9922	32	1.0383	-1.0361
15	1.5050	-1.5281	32	0.5853	0.4341
15	0.5091	0.9703	33	-0.1648	0.2045
16	1.0020	-0.9355	33	0.4332	-0.5981
16	1.4110	-0.4322	34	-0.9889	0.9534
17	0.2890	-0.3170	34	-0.9984	0.0301
17	-0.2156	0.5015	35	0.9476	-0.9653
18	-1.5652	1.5185	35	0.2739	0.6579

After exclusion of our suspect points from calculation of parameters in the “Crystal ball” model, we obtain the results shown in the table 7.

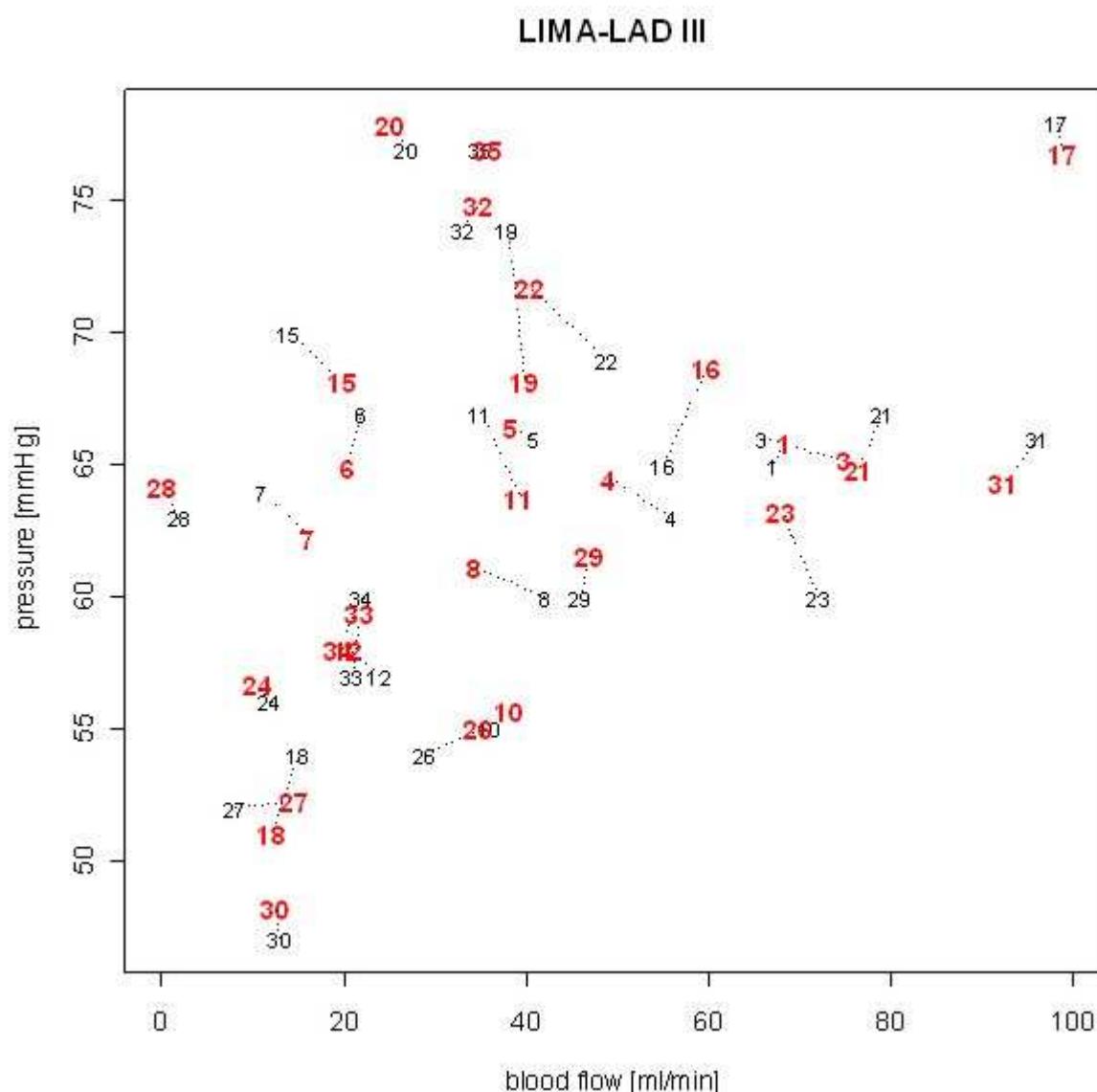
When we compare values of calculated parameters in the case of an entire data set and a data set with exclusion of suspect outliers, we can see that the blood flow after removal of the cross clamp depends more on the blood flow at the time when the cross clamp is applied to the aorta and much less on the blood pressure at the time when the cross clamp is in place on the aorta. Moreover, we can claim that the differences in dependence on blood flow (less in

A Prediction of Blood Flow through a Bypass Graft Using Statistical Methods

the entire data set) and blood pressure (higher in the entire data) get bigger in a case of blood pressure after removal of the cross clamp.

Table 7. Estimated parameters after exclusion of the suspect outliers.

$\hat{\beta}_1$	-45.67
$\hat{\beta}_2$	-19.17
$\hat{\beta}_3$	1.04
$\hat{\beta}_4$	0.50
$\hat{\beta}_5$	0.05
$\hat{\beta}_6$	1.24



*Fig. 3. Predicted and measured values of blood pressure and blood flow through a bypass after removal of the cross clamp from the aorta after exclusion of suspect outliers.*

Already in the figure 3, where estimates of blood flow are marked in bold red and blood pressure and in black, we can see that the residuals got lower after the exclusion of suspect outliers. Also, the table of residuals confirms this statement.

*Table 8. Residuals after exclusion of the suspect outliers.*

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
1	-1.40	-0.86
3	-8.74	0.83
4	6.91	-1.50
5	2.66	-0.45
6	1.55	2.04
7	-5.27	1.73
8	7.68	-1.15
10	-2.13	-0.67
11	-4.31	3.19
12	3.48	-0.98

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
15	-5.96	1.80
16	-4.77	-3.69
17	-0.88	1.15
18	2.84	2.99
19	-1.79	5.79
20	2.00	-0.97
21	2.60	2.13
22	8.55	-2.76
23	4.09	-3.29
24	1.35	-0.64

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
26	-5.90	-0.99
27	-6.69	-0.22
28	1.94	-1.19
29	-0.87	-1.60
30	0.54	-1.13
31	3.57	1.65
32	-1.90	-0.86
33	-0.70	-2.34
34	2.42	1.95
35	-0.86	0.04

## 2.4 Location of leverage points in the “Crystal ball” model

In this section, we find points which most affect the values of estimated blood flow and blood pressure in the “Crystal ball” model (section 1.2, equations (2), (3)).

We want to find the maximum of absolute values of derivation of the terms (38) and (39), i.e. maximum of:

$$\left| \frac{\partial(\Pi_1 Z)_i}{\partial(Z)_j} \right| = |(\Pi_1)_{i,j}| \quad \text{and} \quad \left| \frac{\partial \hat{\beta}_i}{\partial(Z)_j} \right| = |(\Pi_2)_{i,j}|. \quad (7)$$

For our data set without the exclusion of suspected outliers, we obtain by calculation as a leverage point for estimates of increments  $\delta\mu$ ,  $\delta\nu$  and also for parameters  $\beta$  the point with index 2. After the exclusion of suspected outliers, we obtain as leverage point for estimates of increments  $\delta\mu$  and  $\delta\nu$  the point with index 17 and for estimate of increments of parameters  $\beta$  the point with index 30.

## Conclusions

If we use the classical linear model for estimate of blood pressure and blood flow after removal of the cross clamp from the aorta during CABG surgery, we obtain less accurate estimates than when we use the “Crystal ball” model. Moreover, when we exclude suspected outliers our estimates are even more accurate. This statement leads to the hypothesis that patients are divided into several groups and with our algorithm we can calculate certain values of parameters for one group and generally different values for other groups. However, inclusion of patients in the wrong group may be manifested as outliers. To confirm this hypothesis we must complete a deeper analysis of a larger group of patients and include in the surveyed parameters those which can affect blood flow through a bypass before and after removal of the cross clamp from the aorta in the way that they define the above mentioned groups. We can use for example the percentage of stenosis in the target coronary bed, the ejection fraction, FFT ratio, the type of bypass used (LIMA, RIMA, SVG), time of dilatation, the bypass as pedicle/scelet and others [2], [3], [12].

In this orientation period of research a very simplified structure of covariance matrix was used. In further research the covariation matrix, which respects differences in dispersions of the measure of blood flow and the measure of blood pressure must be studied.

The numerical results obtained by suggested method the “Crystal ball” indicate real possibility to predict values of blood flow and blood pressure after removal of the cross clamp from the aorta. However, many calculations on larger data sets must be done before release of this method.

## Appendix The “Crystal ball” model

The model of the relationship between blood flow and pressure at the time when the clamp is applied to the aorta (LIMA-LAD I) and at the time when the clamp is removed (LIMA-LAD III) we assume in this form (for notation see section 1.1):

$$E(\zeta) = [\square_{1,1}, \square_{1,2}, \dots, \square_{n,1}, \square_{n,2}, \square_{1,1}, \square_{1,2}, \dots, \square_{n,1}, \square_{n,2}], \quad Var(\zeta) = \begin{pmatrix} 2n\Theta & \mathbf{I} \\ \mathbf{I} & 2n,2n \end{pmatrix}, \quad (8)$$

where  $E(\zeta)$  satisfies the conditions:

$$\begin{pmatrix} E(\xi_i) \\ E(\eta_i) \end{pmatrix} = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} \beta_3 & \beta_4 \\ \beta_5 & \beta_6 \end{pmatrix} \begin{pmatrix} E(x_i) \\ E(y_i) \end{pmatrix}, \quad i=1,2,\dots,n. \quad (9)$$

Note:  $E(\zeta)$  and  $Var(\zeta)$  denote the mean value and the variance of the random vector  $\zeta$ .

## A.1 Estimate of parameters of the model

The task is to estimate values of model parameters, i.e.  $\mu_{1,1}, \mu_{1,2}, \dots, \mu_{n,1}, \mu_{n,2}$ ,  $v_{1,1}, v_{1,2}, \dots, v_{n,1}, v_{n,2}$  and namely  $\beta_1, \dots, \beta_6$ , in terms of our measured data. To solve this problem we use linearization of the model. The model is not linear because there are products of parameters  $\beta$  and  $\mu$ . After linearization we use estimate algorithms from the model of a direct incomplete measurement of a vector parameter with a system of constraints from the fourth section of the book Statistika a metrologie [14] and from the book Statistical models with Linear Structures [15]. A necessary condition for estimation of the parameters is the number of observations plus the number of constraints has to be greater than a number of the parameters, in our case:  $4n+2n>4n+6$ , because we measure four parameters of each patient.

After linearization we get the model of a direct incomplete measurement of a vector parameter:

$$\zeta \square \begin{pmatrix} \square_{1,1}^{(0)} \\ \square_{1,2}^{(0)} \\ \vdots \\ \square_{n,1}^{(0)} \\ \square_{n,2}^{(0)} \end{pmatrix} \sim_{140} \left( \begin{pmatrix} \square_{1,1} \\ \square_{1,2} \\ \vdots \\ \square_{n,1} \\ \square_{n,2} \end{pmatrix}; \square^2 \mathbf{I} \right) \quad (10)$$

with constraints

$$\mathbf{b} = \mathbf{B}_1 \begin{pmatrix} \square_{1,1} \\ \square_{1,2} \\ \vdots \\ \square_{n,1} \\ \square_{n,2} \end{pmatrix} \quad \mathbf{B}_2 \begin{pmatrix} \square_1 \\ \vdots \\ \square_6 \end{pmatrix} = 0 \quad (11)$$

A Prediction of Blood Flow through a Bypass Graft Using Statistical Methods

where

$$\mathbf{b} = \begin{pmatrix} -\nu_{1,1}^{(0)} + \beta_1^{(0)} + \beta_3^{(0)} \mu_{1,1}^{(0)} + \beta_4^{(0)} \mu_{1,2}^{(0)} \\ \vdots \\ -\nu_{i,1}^{(0)} + \beta_1^{(0)} + \beta_3^{(0)} \mu_{i,1}^{(0)} + \beta_4^{(0)} \mu_{i,2}^{(0)} \\ -\nu_{i,2}^{(0)} + \beta_2^{(0)} + \beta_4^{(0)} \mu_{i,1}^{(0)} + \beta_6^{(0)} \mu_{i,2}^{(0)} \\ \vdots \\ -\nu_{i,2}^{(0)} + \beta_2^{(0)} + \beta_4^{(0)} \mu_{i,1}^{(0)} + \beta_6^{(0)} \mu_{i,2}^{(0)} \end{pmatrix}. \quad (12)$$

The matrix  $\mathbf{B}_1$  of the type  $2n \times 4n$  is (a vertical line separates first  $2n$  columns):

$$\left( \begin{array}{cccc|cccc|ccccccccc} \beta_3^{(0)} & \beta_4^{(0)} & 0 & 0 & \cdots & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & \cdots & 0 & 0 & 0 \\ \beta_5^{(0)} & \beta_6^{(0)} & 0 & 0 & \cdots & 0 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & \cdots & 0 & 0 & 0 \\ 0 & 0 & \beta_3^{(0)} & \beta_4^{(0)} & 0 & \cdots & 0 & 0 & 0 & 0 & -1 & 0 & 0 & \cdots & 0 & 0 & 0 \\ 0 & 0 & \beta_5^{(0)} & \beta_6^{(0)} & 0 & \cdots & 0 & 0 & 0 & 0 & 0 & -1 & 0 & \cdots & 0 & 0 & 0 \\ \cdots & \cdots \\ 0 & 0 & 0 & 0 & 0 & \cdots & 0 & \beta_3^{(0)} & \beta_4^{(0)} & 0 & 0 & 0 & 0 & \cdots & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & \cdots & 0 & \beta_5^{(0)} & \beta_6^{(0)} & 0 & 0 & 0 & 0 & \cdots & 0 & 0 & -1 \end{array} \right). \quad (13)$$

The matrix  $\mathbf{B}_2$  of the type  $2n \times 6$  is:

$$\left( \begin{array}{cccccc} 1 & 0 & \square_{1,1}^{(0)} & \square_{1,2}^{(0)} & 0 & 0 \\ 0 & 1 & 0 & 0 & \square_{1,1}^{(0)} & \square_{1,2}^{(0)} \\ 1 & 0 & \square_{2,1}^{(0)} & \square_{2,2}^{(0)} & 0 & 0 \\ 0 & 1 & 0 & 0 & \square_{2,1}^{(0)} & \square_{2,2}^{(0)} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & 0 & \square_{n,1}^{(0)} & \square_{n,2}^{(0)} & 0 & 0 \\ 0 & 1 & 0 & 0 & \square_{n,1}^{(0)} & \square_{n,2}^{(0)} \end{array} \right) \quad (14)$$

Estimations of parameters  $\delta\mu_{i,j}$ ,  $\delta\nu_{i,j}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  and  $\delta\beta_i$ ,  $i=1,2,\dots,6$  are obtained by minimization of function:

$$\boxed{\square \boxed{\square}_{1,1}, \dots, \boxed{\square}_{n,2}} \left( \zeta \boxed{\begin{pmatrix} \boxed{\square}_{1,1}^{(0)} \\ \vdots \\ \boxed{\square}_{n,2}^{(0)} \end{pmatrix}} \right) \boxed{\begin{pmatrix} \boxed{\square}_{1,1} \\ \vdots \\ \boxed{\square}_{n,2} \end{pmatrix}} \left( \zeta \boxed{\begin{pmatrix} \boxed{\square}_{1,1}^{(0)} \\ \vdots \\ \boxed{\square}_{n,2}^{(0)} \end{pmatrix}} \right) \boxed{\begin{pmatrix} \boxed{\square}_{1,1} \\ \vdots \\ \boxed{\square}_{n,2} \end{pmatrix}} \quad (15)$$

with satisfaction of the constraints (11).

After tedious calculation we get the result in the form:

$$\begin{pmatrix} \boxed{\hat{\square}}_{1,1} \\ \boxed{\hat{\square}}_{1,2} \\ \vdots \\ \boxed{\hat{\square}}_{n,1} \\ \boxed{\hat{\square}}_{n,2} \\ \boxed{\hat{\square}}_{1,1} \\ \boxed{\hat{\square}}_{1,2} \\ \vdots \\ \boxed{\hat{\square}}_{n,1} \\ \boxed{\hat{\square}}_{n,2} \end{pmatrix} = \mathbf{Z} \boxed{\mathbf{B}_1 \boxed{\mathbf{B}_1 \mathbf{B}_1}^T \mathbf{B}_2 \mathbf{B}_2^T \mathbf{B}_1 \mathbf{B}_1^T \mathbf{B}_2 \mathbf{B}_2^T \mathbf{B}_1 \mathbf{B}_1^T \mathbf{B}_2 \mathbf{B}_2^T \mathbf{B}_2 \mathbf{B}_2^T \mathbf{B}_1 \mathbf{B}_1^T \mathbf{B}_2 \mathbf{B}_2^T \mathbf{B}_1 \mathbf{Z}}, \quad (16)$$

$$(\hat{\delta\beta}_1, \hat{\delta\beta}_2, \hat{\delta\beta}_3, \hat{\delta\beta}_4, \hat{\delta\beta}_5, \hat{\delta\beta}_6)' = -[\mathbf{B}_2' (\mathbf{B}_1 \mathbf{B}_1' + \mathbf{B}_2 \mathbf{B}_2')^{-1} \mathbf{B}_2]^{-1} \mathbf{B}_2' (\mathbf{B}_1 \mathbf{B}_1' + \mathbf{B}_2 \mathbf{B}_2')^{-1} \mathbf{B}_1 \mathbf{Z}, \quad (17)$$

where

$$\mathbf{Z} \boxed{[x_1 \boxed{\square}_{1,1}^{(0)}, y_1 \boxed{\square}_{1,2}^{(0)}, \dots, x_n \boxed{\square}_{n,1}^{(0)}, y_n \boxed{\square}_{n,2}^{(0)}, \boxed{\square}_1 \boxed{\square}_{1,1}^{(0)}, \boxed{\square}_1 \boxed{\square}_{1,2}^{(0)}, \dots, \boxed{\square}_n \boxed{\square}_{n,1}^{(0)}, \boxed{\square}_n \boxed{\square}_{n,2}^{(0)}]} \quad (18)$$

This result is considered as the result of the first iteration step, i.e.:

$$\begin{pmatrix} \hat{\square}_{1,1}^{(1)} \\ \vdots \\ \hat{\square}_{n,2}^{(1)} \end{pmatrix} \begin{pmatrix} \boxed{\square}_{1,1}^{(0)} \\ \vdots \\ \boxed{\square}_{n,2}^{(0)} \end{pmatrix} \begin{pmatrix} \boxed{\hat{\square}}_{1,1} \\ \vdots \\ \boxed{\hat{\square}}_{n,2} \end{pmatrix}; \begin{pmatrix} \hat{\square}_1^{(1)} \\ \vdots \\ \hat{\square}_6^{(1)} \end{pmatrix} \begin{pmatrix} \boxed{\square}_1^{(0)} \\ \vdots \\ \boxed{\square}_6^{(0)} \end{pmatrix} \begin{pmatrix} \boxed{\hat{\square}}_1 \\ \vdots \\ \boxed{\hat{\square}}_6 \end{pmatrix} \quad (19)$$

Following vectors are used in the next iteration step

$$\begin{pmatrix} \hat{\mu}_{1,1}^{(1)} \\ \vdots \\ \hat{\nu}_{n,2,k\text{orig}}^{(1)} \end{pmatrix}, \begin{pmatrix} \hat{\beta}_1^{(1)} \\ \vdots \\ \hat{\beta}_6^{(1)} \end{pmatrix} \text{ instead of vectors } \begin{pmatrix} \mu_{1,1}^{(0)} \\ \vdots \\ \nu_{n,2}^{(0)} \end{pmatrix}, \begin{pmatrix} \beta_1^{(0)} \\ \vdots \\ \beta_6^{(0)} \end{pmatrix}.$$

The resulting estimate can be written as follows:

$$\begin{array}{ll} \hat{\square}_{i,1}^{(k)} & \hat{\square}_{i,1}^{(k\Box 1)} & \hat{\square}_{i,1}^{(k)} \\ \hat{\square}_{i,2}^{(k)} & \hat{\square}_{i,2}^{(k\Box 1)} & \hat{\square}_{i,2}^{(k)} \end{array} \text{ and } \begin{array}{ll} \hat{\square}_{i,1}^{(k)} & \hat{\square}_{i,1}^{(k\Box 1)} & \hat{\square}_{i,1}^{(k)} \\ \hat{\square}_{i,2}^{(k)} & \hat{\square}_{i,2}^{(k\Box 1)} & \hat{\square}_{i,2}^{(k)} \end{array} i=1,2,\dots,n, \quad (20)$$

$$\hat{\square}_i^{(k)} \quad \hat{\square}_i^{(k\Box 1)} \quad \hat{\square}_i^{(k)} \quad i=1,2,\dots,6,$$

where  $\hat{\square}_{i,j}^{(k\Box 1)}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  and  $\hat{\square}_{i,j}^{(k\Box 1)}$ ,  $i=1,2,\dots,6$  are estimates of parameters  $\nu$  and  $\beta$  from the previous iteration.

The vector  $\mathbf{Z}^{(k)}$  in the k-th iteration step is:

$$[\mathbf{Z}^{(k)}]^T = [x_1 \square \hat{\square}_{1,1}^{(k\Box 1)}, y_1 \square \hat{\square}_{1,2}^{(k\Box 1)}, \dots, x_n \square \hat{\square}_{n,1}^{(k\Box 1)}, y_n \square \hat{\square}_{n,2}^{(k\Box 1)}, \hat{\square}_1 \square \hat{\square}_{1,1}^{(k\Box 1)}, \hat{\square}_1 \square \hat{\square}_{1,2}^{(k\Box 1)}, \dots, \hat{\square}_n \square \hat{\square}_{n,1}^{(k\Box 1)}, \hat{\square}_n \square \hat{\square}_{n,2}^{(k\Box 1)}] \quad (21)$$

In the iterative process we correct the estimates of the parameters  $\nu$  with these relationships:

$$\begin{aligned} \hat{\nu}_{i,1,k\text{orig}}^{(k)} &= \hat{\beta}_1^{(k)} + \hat{\beta}_3^{(k)} \hat{\mu}_{i,1}^{(k)} + \hat{\beta}_4^{(k)} \hat{\mu}_{i,2}^{(k)} & i=1,2,\dots,n, k=0,1,2\dots \\ \hat{\nu}_{i,2,k\text{orig}}^{(k)} &= \hat{\beta}_2^{(k)} + \hat{\beta}_5^{(k)} \hat{\mu}_{i,1}^{(k)} + \hat{\beta}_6^{(k)} \hat{\mu}_{i,2}^{(k)} \end{aligned} \quad (22)$$

For the initial (zero) iteration we choose

$$\hat{\square}_{i,1}^{(0)} = x_i, \quad \hat{\square}_{i,2}^{(0)} = y_i \quad \text{for } i=1,2,\dots,n \quad (23)$$

and parameters  $\nu_{i,j}^{(0)}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  and  $\beta_i^{(0)}$ ,  $i=1,2,\dots,6$  we calculate from equations for three points (in the xy plane), which lay at borders of points field, i.e. for certain coordinates:

*A Prediction of Blood Flow through a Bypass Graft Using Statistical Methods*

$\begin{pmatrix} x_{i1} \\ y_{i1} \end{pmatrix}, \begin{pmatrix} x_{i2} \\ y_{i2} \end{pmatrix}, \begin{pmatrix} x_{i3} \\ y_{i3} \end{pmatrix}$  and corresponding points  $\begin{pmatrix} \xi_{i1} \\ \eta_{i1} \end{pmatrix}, \begin{pmatrix} \xi_{i2} \\ \eta_{i2} \end{pmatrix}, \begin{pmatrix} \xi_{i3} \\ \eta_{i3} \end{pmatrix}$ .

For these points we solve the system of equations:

$$\begin{pmatrix} 1 & 0 & x_{i1} & y_{i1} & 0 & 0 \\ 0 & 1 & 0 & 0 & x_{i1} & y_{i1} \\ 1 & 0 & x_{i2} & y_{i2} & 0 & 0 \\ 0 & 1 & 0 & 0 & x_{i2} & y_{i2} \\ 1 & 0 & x_{i3} & y_{i3} & 0 & 0 \\ 0 & 1 & 0 & 0 & x_{i3} & y_{i3} \end{pmatrix} \begin{pmatrix} \beta_1^{(0)} \\ \beta_2^{(0)} \\ \beta_3^{(0)} \\ \beta_4^{(0)} \\ \beta_5^{(0)} \\ \beta_6^{(0)} \end{pmatrix} = \begin{pmatrix} \xi_{i1} \\ \eta_{i1} \\ \xi_{i2} \\ \eta_{i2} \\ \xi_{i3} \\ \eta_{i3} \end{pmatrix}. \quad (24)$$

Parameters  $\beta_i^{(0)}$ ,  $i=1,2,\dots,6$  are the solution of this system and they are used for the calculation of the initial iteration of parameters  $v_{i,j}^{(0)}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  in following equations:

$$\begin{matrix} \square_{i,1}^{(0)} & \square_1^{(0)} & \square_3^{(0)} & \square_{i,1}^{(0)} & \square_4^{(0)} & \square_{i,2}^{(0)} \\ \square_{i,2}^{(0)} & \square_2^{(0)} & \square_5^{(0)} & \square_{i,1}^{(0)} & \square_6^{(0)} & \square_{i,2}^{(0)} \end{matrix} \quad i=1,2,\dots,n, \quad (25)$$

where parameters  $\mu_{i,j}^{(0)}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  are set up in this way:

$$\square_{i,1}^{(0)} = x_i \quad \text{and} \quad \square_{i,2}^{(0)} = y_i. \quad (26)$$

The iterative calculation is finished when estimates of the increments  $\delta\mu_{i,j}$ ,  $\delta v_{i,j}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  and  $\delta\beta_i$ ,  $i=1,2,\dots,6$  can be omitted, i.e. on the following condition:

$$|\hat{\mu}_{i,j} - \hat{\mu}_{i,j,korig}| \leq \varepsilon \quad i=1,2,\dots,n, j=1,2 \quad (27)$$

where  $\varepsilon$  is a preset constant.

Estimates of parameters  $\beta_i$  obtained in six iterative steps ( $\varepsilon=0,1$ ) and our data set are shown in the table 4.

An estimator of a covariance matrix for these parameters is calculated by following relationship:

*A Prediction of Blood Flow through a Bypass Graft Using Statistical Methods*

$$\hat{Var} \begin{pmatrix} \hat{\square}_1 \\ \hat{\square}_2 \\ \hat{\square}_3 \\ \hat{\square}_4 \\ \hat{\square}_5 \\ \hat{\square}_6 \end{pmatrix} = \hat{\sigma}^2 \left[ \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{I} \right] \quad (28)$$

where

$$\hat{\sigma}^2 = \frac{\mathbf{v}_1^\top \mathbf{v}_1 + \mathbf{v}_2^\top \mathbf{v}_2}{64} \quad (29)$$

whereas

64 = number of measurements(140)+number of constraints(70)–number of parameters(140+6)

and

$$\begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{pmatrix} = \begin{pmatrix} x_1 & \hat{\square}_{1,1} \\ \vdots & \vdots \\ y_n & \hat{\square}_{n,2} \\ \hat{\square}_1 & \hat{\square}_{1,1} \\ \vdots & \vdots \\ \hat{\square}_n & \hat{\square}_{n,2} \end{pmatrix} \quad (30)$$

For our data set we obtain following numerical results:

$$\hat{\sigma}^2 = 46.52034,$$

$$\hat{Var} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \\ \hat{\beta}_4 \\ \hat{\beta}_5 \\ \hat{\beta}_6 \end{pmatrix} = \begin{pmatrix} 539.914 & 225.030 & -0.840 & -7.528 & -0.350 & -3.137 \\ 225.030 & 326.432 & -0.350 & -3.137 & -0.508 & -4.551 \\ -0.840 & -0.350 & 0.008 & 0.005 & 0.003 & 0.002 \\ -7.528 & -3.137 & 0.005 & 0.110 & 0.002 & 0.046 \\ -0.350 & -0.508 & 0.003 & 0.002 & 0.005 & 0.003 \\ -3.137 & -4.551 & 0.002 & 0.046 & 0.003 & 0.066 \end{pmatrix}.$$

## A.2 Location of outlier points in the “Crystal ball” model

When the iteration is stopped (k-th step), for our vector of residuals  $\mathbf{Z}^{(k)}$  from the equation (21), the following relationship is certainly valid :

$$\mathbf{Z}^{(k)} \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{pmatrix} \begin{pmatrix} x_1 & \square \hat{\Delta}_{1,1}^{(k)} \\ \vdots & \vdots \\ y_n & \square \hat{\Delta}_{n,2}^{(k)} \\ \hline \square_1 & \square \hat{\Delta}_{1,1}^{(k)} \\ \vdots & \vdots \\ \square_n & \square \hat{\Delta}_{n,2}^{(k)} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ y_n \\ \hline \square_1 \\ \vdots \\ \square_n \end{pmatrix} \square \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \hat{\mu}^{(k)} \\ \hat{\mathbf{v}}^{(k)} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ y_n \\ \hline \square_1 \\ \vdots \\ \square_n \end{pmatrix} \square \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mu^{(k-1)} & \delta\hat{\mu}^{(k)} \\ \mathbf{v}^{(k-1)} & \delta\hat{\mathbf{v}}^{(k)} \end{pmatrix}$$

$$\begin{pmatrix} x_1 \\ \vdots \\ y_n \\ \hline \square_1 \\ \vdots \\ \square_n \end{pmatrix} \square \mu^{(k-1)} \square \begin{pmatrix} x_1 \\ \vdots \\ y_n \\ \hline \square_1 \\ \vdots \\ \square_n \end{pmatrix} \square \mathbf{v}^{(k-1)} \quad (31)$$

$$\square \left[ \mathbf{Z}^{(k-1)} \square \mathbf{B}_1^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_1 \mathbf{Z}^{(k-1)} \right]$$

When we substitute the vector  $\mathbf{Z}$  from the last iteration into the last row of expression we can rewrite the relationship in the following form:

$$\begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{pmatrix} \square \mathbf{B}_1^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_1 \mathbf{Z}^{(k-1)} \quad (32)$$

Thereafter, when we denote

$$\mathbf{T} = \left[ \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \right] \quad (33)$$

we can write

$$Var \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{pmatrix} = \square^2 \mathbf{B}_1^\top \mathbf{T} \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{T} \mathbf{B}_1 = \begin{pmatrix} Var(\mathbf{v}_1) & cov(\mathbf{v}_1, \mathbf{v}_2) \\ cov(\mathbf{v}_2, \mathbf{v}_1) & Var(\mathbf{v}_2) \end{pmatrix} = \square^2 \begin{pmatrix} \mathbf{Q}_{11} & \mathbf{Q}_{12} \\ \mathbf{Q}_{21} & \mathbf{Q}_{22} \end{pmatrix} \quad (34)$$

For every  $i=1,2,\dots,2n$  we calculate these expressions:

$$|\mathbf{v}_1| / \sqrt{Var(\mathbf{v}_1)}_i \quad |\mathbf{v}_1| / \sqrt{\square^2 \mathbf{Q}_{11}}_i \quad (35)$$

$$|\mathbf{v}_2| / \sqrt{Var(\mathbf{v}_2)}_i \quad |\mathbf{v}_2| / \sqrt{\square^2 \mathbf{Q}_{22}}_i$$

### A.3 Location of leverage points in the “Crystal ball” model

When the iteration process is stopped (k-th step) we use equations for increments

$$\begin{pmatrix} \hat{\alpha}_{1,1} \\ \hat{\alpha}_{1,2} \\ \vdots \\ \hat{\alpha}_{n,1} \\ \hat{\alpha}_{n,2} \\ \hat{\alpha}_{1,1} \\ \hat{\alpha}_{1,2} \\ \vdots \\ \hat{\alpha}_{n,1} \\ \hat{\alpha}_{n,2} \end{pmatrix} = \mathbf{Z}^{(k)} - \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{Z}^{(k)} \quad (36)$$

$$[\hat{\alpha}_1, \hat{\alpha}_2, \hat{\alpha}_3, \hat{\alpha}_4, \hat{\alpha}_5, \hat{\alpha}_6]^\top = \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{Z}^{(k)} \quad (37)$$

which we can write as follows:

$$\Pi_1 \mathbf{Z}^{(k)} = \begin{pmatrix} \hat{\alpha}_{1,1} \\ \hat{\alpha}_{1,2} \\ \vdots \\ \hat{\alpha}_{n,1} \\ \hat{\alpha}_{n,2} \\ \hat{\alpha}_{1,1} \\ \hat{\alpha}_{1,2} \\ \vdots \\ \hat{\alpha}_{n,1} \\ \hat{\alpha}_{n,2} \end{pmatrix} = \mathbf{I} - \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{Z}^{(k)} \quad (38)$$

$$\Pi_2 \mathbf{Z}^{(k)} = [\hat{\alpha}_1, \hat{\alpha}_2, \hat{\alpha}_3, \hat{\alpha}_4, \hat{\alpha}_5, \hat{\alpha}_6]^\top = \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{Z}^{(k)} \quad (39)$$

## Bibliography

- [1] ÚZIS: Aktuální informace č. 14 - Kardiochirurgické operace ([http://www.uzis.cz/cz/archiv04/14\\_04.pdf](http://www.uzis.cz/cz/archiv04/14_04.pdf))
- [2] Takami, Y., Ina, H.: A simple Method to Determine Anastomotic Quality of Coronary Artery Bypass Grafting in the Operating Room. *Cardiovascular Surgery*, Vol. 9, No. 5, pp. 499-503
- [3] Takami, Y., Ina, H.: Relation of Intraoperative Flow Measurement With Postoperative Quantitative Angiographic Assessment of Coronary Artery Bypass Grafting. *Ann Thorac Surg* 2001, 72:1270-4
- [4] D'Anconna, G., Karamanoukian, H. L., Bergsland, J.: Is Intraoperative measurement of coronary blood flow a good predictor of graft patency? *European Journal of Cardio-thoracic Surgery* 20 (2001), pp. 1075-1076
- [5] Shin, H. et al.: Intraoperative Assessment of Coronary Artery Bypass Graft: Transit-Time Flowmetry Versus Angiography. *Ann Thorac Surg* 2001, 72:1562-5
- [6] Louagie, Y. et al.: Pulsed Doppler Intraoperative Flow Assessment and Midterm Coronary Graft Patency. *Ann Thorac Surg* 1998, 66:1282-8
- [7] Louagie, Y., Brockmann, C., Gurné, O., Jamart, J.: Intraoperative Flow Measurements: Predictive Value for Postoperative Angiographic Follow-Up. In.: *Intraoperative Graft Patency Verification in Cardiac and Vascular Surgery* (ed. G. D'Ancona). Futura Publishing, Armonk (NY) 2001
- [8] Milnor, W. R.: *Hemodynamics*. Williams&Wilkins, Baltimore (USA) 1998
- [9] Hata, M. et al.: Midterm Results of Coronary Artery Bypass Graft Surgery With Internal Thoracic Artery Under Low Free-Flow Conditions. *Ann Thorac Surg* 2004, 78:477-80
- [10] Pagni, S. et al.: Factors Affecting Internal Mammary Artery Graft Survival: How Is Competitive Flow from a Patent Native Coronary Vessel a Risk Factor? *Journal of Surgical Research* 71 (1997), pp. 172-178
- [11] Sabik, J. F. et al.: Does Competitive Flow Reduce Internal Thoracic Artery Graft Patency? *Ann Thorac Surg* 2003, 76:1490-7
- [12] Pagni, S., Storey, J. et al.: ITA versus SVG: a comparison of instantaneous pressure and flow dynamics during competitive flow. *European Journal of Cardio-thoracic Surgery* 11 (1997), pp. 1086-1092
- [13] Lausten, J.: *Transit-Time Flow Measurement: Principles and Clinical Applications*. In.: *Intraoperative Graft Patency Verification in Cardiac and Vascular Surgery* (ed. G. D'Ancona). Futura Publishing, Armonk (NY) 2001
- [14] Kubáček, L., Kubáčková, L.: *Statistika a metrologie*. VUP Olomouc 2000
- [15] Kubáček, L., Kubáčková, L., Volaufová, J.: *Statistical Models with Linear Structures*. Veda, Bratislava 1995.

# Předpověď průtoku krve bypassem pomocí statistických metod

Jana Vrbková<sup>1</sup>, Vilém Bruk<sup>2</sup>

*1. Department of Mathematical Analysis and Applications of Mathematics, Faculty of Natural Sciences, Palacky University Olomouc, Czech Republic, 2. Clinic of Cardiosurgery, Faculty Hospital Olomouc, Czech Republic*

Revaskulatizace myokardu patří mezi nejčastější kardiochirurgické zákroky. Perioperační i dlouhodobá úspěšnost revaskularizace koronárního řečiště závisí na průchodnosti použitého štěpu a kvalitě anastomózy. Hemodynamické charakteristiky měřené peroperačně pomáhají ověřit kvalitu anastomózy a ovlivňují i dlouhodobou průchodnost použitého štěpu. Při operacích bypassu v mimotělním oběhu chirurg měří průtok krve štěpem v době, kdy je na aortě naložená svorka (do srdce nepřichází krev prostřednictvím koronárních cév, ale pouze měřeným štěpem), a později měří na tomtéž místě průtok po povolení svorky na aortě. Cílem této práce je nalézt model s jehož pomocí bude možné předpovědět průtok krve štěpem po povolení svorky na aortě na základě prvního měření při naložené svorce. Při spolehlivé predikci by bylo možné snížit počet měření se současným zachováním informace o objektu.

**Klíčová slova:** revaskularizace myokardu, predikce průtoku a tlaku krve, vícerozměrná lineární regrese, nelineární regresní model, linearizace, lineární regresní model s podmínkami, outlier, leverage

## Úvod

Operace bypassu patří v celém světe i v naší republice mezi nejčastější kardiochirurgické zákroky (10 797 srdečních operací, z toho 7051 operací aortokoronárních bypassů v roce 2002 v ČR) [1]. Perioperační i dlouhodobá úspěšnost revaskularizace koronárního řečiště přitom závisí na průchodnosti použitého štěpu a kvalitě anastomózy. Řada autorů se zabývala stanovením charakteristik, které jsou určující pro stanovení průchodnosti štěpu a kvality anastomózy, jak z krátkodobého tak z dlouhodobého hlediska. Ukazuje se, že hemodynamické charakteristiky pomáhají ověřit kvalitu anastomózy a ovlivňují i dlouhodobou průchodnost použitého štěpu [2], [3], [4], [5], [6]. Louagie uvádí jako dominantní charakteristiku ovlivňující dlouhodobou průchodnost štěpu rezistenci, kterou lze stanovit jako podíl tlaku krve a průtoku krve cévou [7]. Tento vztah je ale zjednodušením reality již proto, že průtok krve štěpem není stacionární a krev není Newtonovská kapalina [8]. Přesto má smysl uvažovat průtok krve štěpem a střední arteriální tlak krve jako veličiny rezistenci determinující. Hata [9] na souboru pacientů, u nichž byl naměřen nízký volný průtok arteriálním štěpem, potvrdil již dříve uvažovaný fakt, že levá mammární artérie (LIMA) je schopna přizpůsobovat svůj průměr v závislosti na potřebách zásobení cílového koronárního řečiště krví. Současně je známo, že průtok štěpem je ovlivněn kompetitivním průtokem nativního koronárního řečiště, přičemž se spekuluje o tom, zda tento kompetitivní průtok může způsobit změnu průchodnosti štěpu a zda jsou všechny štěpy (LIMA, RIMA, SVG) stejně citlivé na jeho působení [10], [11], [12]. V průběhu revaskularizace myokardu –

## Předpověď průtoku krve bypassem pomocí statistických metod

operace v mimotělním oběhu, se proto provádí více měření: free-flow (volný průtok štěpem, který ještě není našit na cílovou koronární artérii), před povolením svorky na aortě (není přítomen kompetitivní průtok nativního koronárního řečiště), po povolení svorky na aortě a na konci operace před uzavřením hrudníku. Někdy se před povolením svorky na aortě provádí měření dvakrát: v době, kdy ostatní štěpy nejsou povoleny, a po povolení všech našítych štěpů (zachycuje vliv kolaterál). V Olomoucké fakultní nemocnici je již od roku 2002 pro měření průtoku krve při operacích aortokoronárního bypassu využíván průtokoměr norské společnosti Medi-Stim, který pracuje na tzv. Transit-Time principu. Ten je založen na stanovení časového rozdílu mezi dobou, kterou urazí ultrazvukový signál z vysílače do přijímače proti směru toku krve, a mezi dobou, kterou urazí opačný signál směřující po směru toku krve [13]. Tento přístroj umožňuje sledovat v průběhu měření přímo při operaci aktuální křivku průtoku krve, velikost průtoku krve, průměrný průtok krve, Pulsatility Index PI=(max. průtok - min. průtok)/průměrný průtok a další veličiny. K přístroji lze připojit dvě sondy pro snímání průtoku krve, další 2 sondy snímající tlak krve a k dispozici jsou rovněž 2 vstupy pro další signály, např. pro EKG. Pro zaznamenaná data lze spočítat další charakteristiky, mimo jiné lze provést rychlou Fourierovu transformaci pro libovolnou zaznamenanou křivku. Na menším souboru pacientů (35) byl při operacích bypassu v mimotělním oběhu (CABG, on-pump) zaznamenán střední arteriální tlak krve v a. radialis a také naměřený průtok krve štěpem (levá mammární artérie na r. interventrikularis anterior - LIMA-LAD) v době, kdy je na aortě naložená svorka (do srdce nepřichází krev prostřednictvím koronárních cév ani jinou cestou, ale pouze měřeným štěpem), a později byl zaznamenán tlak krve a změřen průtok na tomtéž místě po povolení svorky na aortě (působení kompetitivního průtoku). Cílem této práce je nalézt model s jehož pomocí bude možné předpovědět průtok krve štěpem po povolení svorky na aortě na základě prvního měření při naložené svorce. Při spolehlivé predikci by bylo možné snížit počet měření se současným zachováním informace o objektu.

## 1 Vstupní data a model

### 1.1 Vstupní data a značení

Pro analýzu jsou použita ručně zaznamenaná data o průtoku krve bypassem tvořeným levou mammární artérií (LIMA) našitou na ramus interventricularis anterior (LAD) levé koronární cévy. Měření bylo prováděno při naložené svorce na aortě (LIMA-LAD I) a při povolené svorce (LIMA-LAD III). Současně s průtokem krve v ml/min byl zaznamenáván i střední arteriální tlak krve v mmHg. Data byla získána měřením u 35 pacientů při operacích LIMA nebo BIMA (měření provedeno na levé mammární artérii).

Data získaná měřením na kardiochirurgické klinice ve FN Olomouc jsou uspořádána do vektoru vstupních dat  $\zeta = (x_1, y_1, \dots, x_n, y_n, \xi_1, \eta_1, \dots, \xi_n, \eta_n)'$

(' označuje transpozici vektoru), kde

$x_1, x_2, \dots, x_n$	značí průtok krve bypassem LIMA-LAD I,
$y_1, y_2, \dots, y_n$	značí střední arteriální tlak krve LIMA-LAD I,
$\xi_1, \xi_2, \dots, \xi_n$	značí průtok krve bypassem LIMA-LAD III,
$\eta_1, \eta_2, \dots, \eta_n$	značí střední arteriální tlak krve LIMA-LAD III.

Předpověď průtoku krve bypassem pomocí statistických metod

Tab. 1. Vstupní (naměřená) data.

i	x <sub>i</sub>	y <sub>i</sub>	$\xi_i$	$\eta_i$
1	80	67	67	65
2	101	47	40	47
3	94	68	66	66
4	53	63	56	63
5	46	66	41	66
6	30	63	22	67
7	34	65	11	64
8	39	60	42	60
9	38	77	58	65
10	55	60	36	55
11	55	63	35	67
12	31	60	24	57
13	85	67	38	65
14	33	80	15	77
15	36	70	14	70
16	74	75	55	65
17	105	72	98	78

18	26	50	15	54
19	51	62	38	74
20	29	77	27	77
21	84	60	79	67
22	40	70	49	69
23	75	65	72	60
24	24	60	12	56
25	44	77	24	75
26	56	62	29	54
27	38	60	8	52
28	10	67	2	63
29	60	65	46	60
30	30	54	13	47
31	99	59	96	66
32	44	76	33	74
33	36	65	21	57
34	31	57	22	60
35	43	76	35	77

## 1.2 Statistický model

Vztah mezi průtokem a tlakem v době naložené svorky na aortě (LIMA-LAD I) a po povolení svorky (LIMA-LAD III) lze předpokládat v několika tvarech. Nejjednodušší možný tvar je klasický lineární regresní model.

$$\begin{pmatrix} v_{i,1} \\ v_{i,2} \end{pmatrix} = \begin{pmatrix} \square_1 \\ \square_2 \end{pmatrix} \begin{pmatrix} \square_{11} & \square_{12} \\ \square_{21} & \square_{22} \end{pmatrix} \begin{pmatrix} x_i \\ y_i \end{pmatrix}, \quad (1)$$

kde  $x_i$  a  $y_i$  jsou hodnoty průtoku a tlaku krve pacienta v čase před povolením svorky na aortě a  $v_{i,1}$  a  $v_{i,2}$  jsou predikované hodnoty těchto parametrů po povolení svorky. Náhodné chyby se v případě použití tohoto modelu vyskytují pouze v hodnotách  $v_{i,1}$  a  $v_{i,2}$ . Parametry  $\alpha_1, \alpha_2, \beta_{11}, \beta_{12}, \beta_{21}, \beta_{22}$  jsou odhadnutý metodu nejmenších čtverců.

Předpoklad o bezchybnosti hodnot  $x_i$  a  $y_i$  však není v souladu s realitou. Model adekvátnější realitě je model „Křišťálová koule“, ve kterém se předpokládá, že i hodnoty  $x_i$  a  $y_i$  jsou realizací náhodných proměnných a tedy:

$$E(\zeta) = [\square_{1,1}, \square_{1,2}, \dots, \square_{n,1}, \square_{n,2}, \square_{1,1}, \square_{1,2}, \dots, \square_{n,1}, \square_{n,2}], \quad Var(\zeta) = \square^2 \begin{pmatrix} 2n, \Theta & I \\ & 2n, 2n \end{pmatrix}, \quad (2)$$

Předpověď průtoku krve bypassem pomocí statistických metod

přičemž  $E(\zeta)$  splňuje podmínky:

$$\begin{pmatrix} E(\square_i) \\ E(\square_i) \end{pmatrix} = \begin{pmatrix} \square_1 \\ \square_2 \end{pmatrix} \begin{pmatrix} \square_3 & \square_4 \\ \square_5 & \square_6 \end{pmatrix} \begin{pmatrix} E(x_i) \\ E(y_i) \end{pmatrix}, \quad i=1,2,\dots,n. \quad (3)$$

Poznámka:  $E(\zeta)$  a  $Var(\zeta)$  značí střední hodnotu a varianci náhodného vektoru  $\zeta$ .

Předpokládaný tvar kovarianční matice  $Var(\zeta)$  je nejjednodušší možný a je použitý kvůli menší obtížnosti dalších výpočtů. Při dalším výzkumu této problematiky bude nutno uvažovat složitější tvar této matice, což však povede k problému odhadu variančních komponent.

## 2 Předpověď hodnot průtoku a tlaku krve

### 2.1 Klasický regresní model

Odhady parametrů  $\alpha_1, \alpha_2, \beta_{11}, \beta_{12}, \beta_{21}, \beta_{22}$  pro naše data získané pomocí funkce `lm(y~x)`, která je dostupná v softwaru R, jsou uvedeny v tab. 2.

Tab. 2. Odhadnuté parametry klasického lineárního modelu.

$\hat{\alpha}$	
[1,]	-37,54
[2,]	5,42

$\hat{\beta}$	[,1]	[,2]
[1,]	0,82	0,044
[2,]	0,52	0,86

Na obr. 1 jsou znázorněny jednak naměřené (černě) ale i odhadnuté (tučně červeně) hodnoty průtoku krve bypassem a tlaku krve po povolení svorky na aortě. Pro srovnání účinnosti regresního modelu a modelu „Křišťálová koule“ jsou v tabulce 3 uvedena rezidua ( $\zeta_i$  a  $\eta_i$  jsou naměřené hodnoty průtoku a tlaku krve po povolení svorky na aortě u i-tého pacienta):

$$\square \square_{i,1} \quad \square_i \square \square_{i,1}, \quad \square \square_{i,2} \quad \square_i \square \square_{i,2}, \quad i=1,2,\dots,n. \quad (4)$$

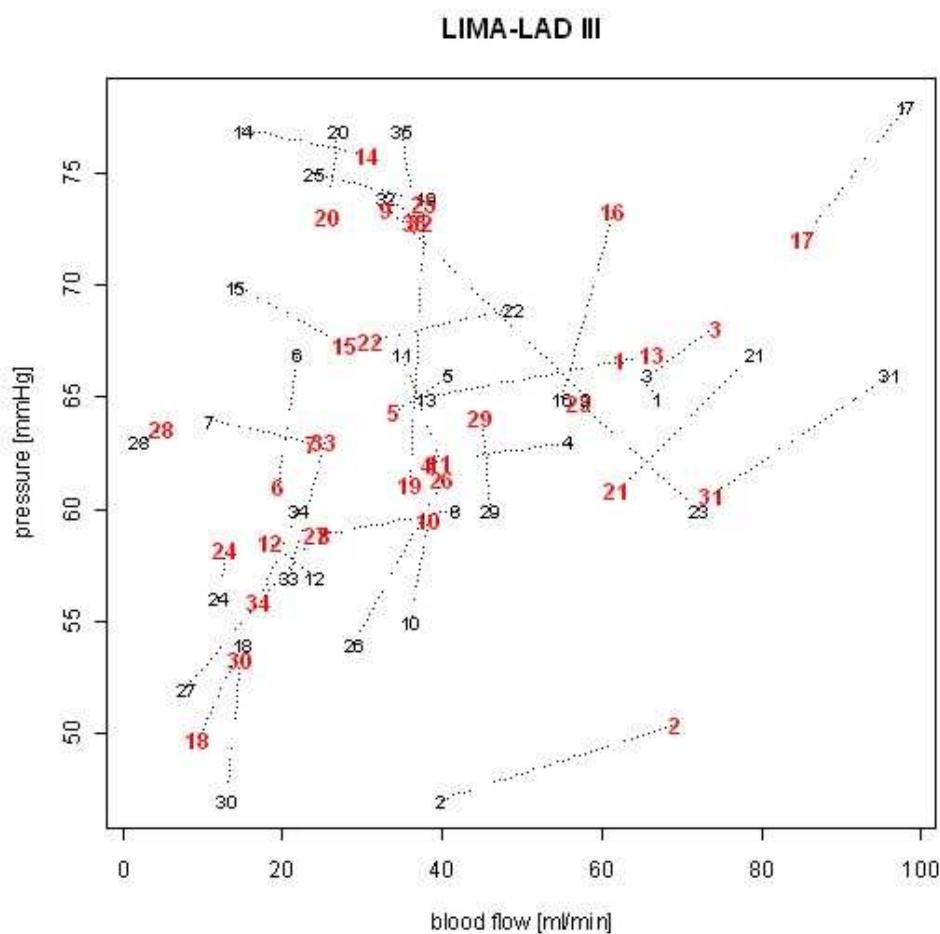
Předpověď průtoku krve bypassem pomocí statistických metod

Tab. 3. Rezidua v klasickém lineárním regresním modelu.

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
1	4,76	-1,69
2	-29,08	-3,36
3	-8,18	-2,16
4	17,85	0,94
5	7,01	1,66
6	2,62	5,95
7	-12,68	1,05
8	16,82	1,14
9	24,88	-8,47
10	-2,24	-4,56
11	-4,78	4,85
12	5,34	-1,51

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
13	-28,32	-1,91
14	-15,59	1,16
15	-13,89	2,65
16	-6,47	-8,32
17	12,78	5,91
18	5,57	4,33
19	1,99	12,89
20	1,22	3,92
21	17,10	6,17
22	17,85	1,47
23	14,87	-4,74
24	-0,94	-2,20

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
25	-14,02	1,26
26	-11,09	-7,33
27	-16,37	-6,82
28	-3,12	-0,63
29	1,11	-4,09
30	-1,75	-6,29
31	22,37	5,38
32	-4,50	1,13
33	-4,31	-6,04
34	4,89	4,08
35	-1,69	4,17



Obr. 1. Předpovězené a naměřené hodnoty tlaku krve a průtoku krve bypassem po povolení svorky na aortě pomocí klasického lineárního regresního modelu.

Předpověď průtoku krve bypassem pomocí statistických metod

## 2.2 Model „Křišťálová koule“

Predikované hodnoty i-tého pacienta určíme ze vztahu

$$\begin{pmatrix} \hat{\beta}_{i,1} \\ \hat{\beta}_{i,2} \end{pmatrix} = \begin{pmatrix} \hat{\beta}_1 & \hat{\beta}_4 \\ \hat{\beta}_2 & \hat{\beta}_5 \end{pmatrix} \begin{pmatrix} \hat{\beta}_3 & \hat{\beta}_4 \\ \hat{\beta}_5 & \hat{\beta}_6 \end{pmatrix} \begin{pmatrix} x_i \\ y_i \end{pmatrix}, \quad (5)$$

kde  $x_i$  a  $y_i$  jsou naměřené hodnoty průtoku a tlaku krve pacienta v čase před povolením svorky na aortě a  $\hat{\beta}_{i,1}$  a  $\hat{\beta}_{i,2}$  jsou predikované hodnoty těchto parametrů po povolení svorky. Postup výpočtu odhadů parametrů  $\beta_i$ ,  $i=1,2,\dots,6$ , který je iterační, je uveden v appendix A.1. Dále bude prověřena účinnost navrhovaného postupu na datech, ze kterých byly odhadnuty parametry  $\beta_i$ . Porovnáme predikované hodnoty s hodnotami přímo naměřenými. Body zvolené pro nultou iteraci ( $i=14, 17, 18$ ) se nacházejí na okrajích pole. Čísla znamenají pořadová čísla případu (pacienta).

Výše uvedenou volbou bodů pro výpočet nulté iterace a konstanty ukončující iterační výpočet  $\epsilon=0,1$  získáme pro naše data po šesti iteračních krocích pro hledané parametry  $\beta_i$  numerické výsledky uvedené v tabulce 4.

Tab. 4. Odhadnuté parametry v modelu „Křišťálová koule“.

$\hat{\beta}_1$	-99,10
$\hat{\beta}_2$	-11,96
$\hat{\beta}_3$	0,96
$\hat{\beta}_4$	1,34
$\hat{\beta}_5$	0,06
$\hat{\beta}_6$	1,12

Pro naše data vypadá kovarianční matice a odhad  $\sigma^2$  takto (postup výpočtu je uveden v appendix A.1):

$$\hat{\sigma}^2 = 46.52034,$$

$$Var \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \\ \hat{\beta}_4 \\ \hat{\beta}_5 \\ \hat{\beta}_6 \end{pmatrix} = \begin{pmatrix} 539.914 & 225.030 & -0.840 & -7.528 & -0.350 & -3.137 \\ 225.030 & 326.432 & -0.350 & -3.137 & -0.508 & -4.551 \\ -0.840 & -0.350 & 0.008 & 0.005 & 0.003 & 0.002 \\ -7.528 & -3.137 & 0.005 & 0.110 & 0.002 & 0.046 \\ -0.350 & -0.508 & 0.003 & 0.002 & 0.005 & 0.003 \\ -3.137 & -4.551 & 0.002 & 0.046 & 0.003 & 0.066 \end{pmatrix}.$$

Předpověď průtoku krve bypassem pomocí statistických metod

Jako numerické vyjádření přesnosti odhadů vypočtených pomocí našeho modelu jsou v tabulce 5 uvedena rezidua:

$$\square \square_{i,1} - \square_i \square \hat{\square}_{i,1}, \quad \square \square_{i,2} - \square_i \square \hat{\square}_{i,2}, \quad i=1,2,\dots,n. \quad (6)$$

Tab. 5. Rezidua v modelu „Křištálová koule“.

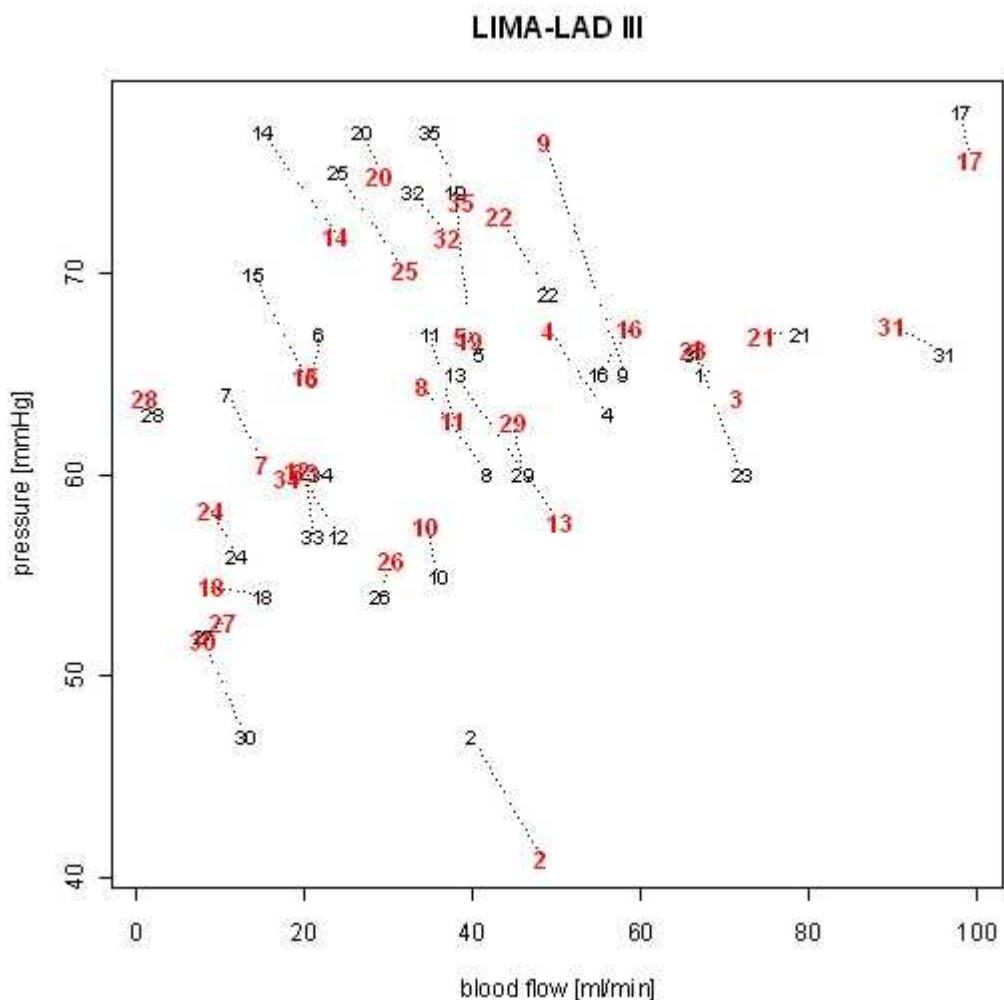
i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
[1]	0,32	-133
[2]	-8,18	6,01
[3]	-5,39	2,25
[4]	6,95	-4,19
[5]	2,33	-0,94
[6]	1,14	2,23
[7]	-4,04	3,42
[8]	7,97	-4,39
[9]	9,40	-11,58
[10]	1,45	-2,46
[11]	-2,67	4,26

[12]	4,77	-3,28
[13]	-12,31	7,37
[14]	-8,76	5,08
[15]	-6,23	5,09
[16]	-3,81	-2,27
[17]	-1,16	2,36
[18]	6,04	-0,41
[19]	-1,75	7,33
[20]	-2,22	2,10
[21]	4,42	0,06
[22]	5,79	-3,86
[23]	5,66	-6,17

[24]	2,96	-2,28
[25]	-8,04	4,90
[26]	-1,69	-1,72
[27]	-2,38	-0,68
[28]	0,76	-0,76
[29]	1,11	-2,61
[30]	4,86	-4,78
[31]	6,06	-1,41
[32]	-4,19	2,26
[33]	0,84	-3,14
[34]	3,89	0,16
[35]	-3,87	3,39

Na obr. 2 jsou znázorněny jednak naměřené (černě) ale i predikované (tučně červeně) hodnoty průtoku krve bypassem a tlaku krve po povolení svorky na aortě.

Předpověď průtoku krve bypassem pomocí statistických metod



Obr. 2. Předpovězené a naměřené hodnoty tlaku krve a průtoku krve bypassem po povolení svorky na aortě v modelu „Křišťálová koule“.

Již z obrázku, ale i hodnot reziduí, je vidět, že model „křišťálová koule“, který má větší počet parametrů, což vede k většímu přiblížení odhadu k odhadovaným hodnotám, je adekvátnější realitě než klasický lineární regresní model.

### 2.3 Nalezení odlehlých („outlier“) bodů v modelu „Křišťálová koule“

V tomto odstavci se budeme snažit hledat odlehlá pozorování v modelu „Křišťálová koule“. Podezřelé body poté z dat použitých pro výpočet parametrů vyloučíme a budeme pozorovat, k jakým změnám v modelu dojde.

Nalezneme takové indexy  $i$ , pro které jsou výrazy (35) větší nebo rovny 1,96 (či přibližně 2). Body s těmito indexy jsou potom naše podezřelé „outliery“.

Pro naše data jsou to body s indexy 2, 9, 13, 14, 25, jak je vidět z tab. 6.

Předpověď průtoku krve bypassem pomocí statistických metod

Tab. 6. Podezřelé „outliery“ v modelu „Křišťálová koule“.

i	$(v_1)_i / \sqrt{[\hat{Var}(v_1)]_0}$	$(v_2)_i / \sqrt{[\hat{Var}(v_2)]_0}$
1	-0,0620	0,0795
1	0,1952	-0,2574
2	<b>2,4664</b>	<b>-2,4933</b>
2	0,9999	1,4239
3	1,3696	-1,3573
3	0,9032	0,4407
4	-1,6880	1,6916
4	-0,8654	-0,7923
5	-0,5717	0,5661
5	-0,3838	-0,1771
6	-0,3219	0,2789
6	-0,7385	0,4257
7	0,9763	-0,9934
7	0,3050	0,6549
8	-1,9367	1,347
8	-1,0731	-0,284
9	<b>-2,2851</b>	<b>2,823</b>
9	0,0311	-2,822
10	-0,3321	0,563
10	0,1412	-0,691
11	0,6068	-0,475
11	-0,2048	0,037
12	-1,1603	1,688
12	-0,5146	-0,253
13	<b>3,490</b>	<b>-3,550</b>
13	1,730	1,212
14	<b>2,011</b>	<b>-2,022</b>
14	1,691	0,922
15	1,050	-1,281
15	0,091	0,703
16	1,020	-0,355
16	1,110	-0,322
17	0,890	-0,170
17	-0,156	0,015
18	-1,652	1,185

i	$(v_1)_i / \sqrt{[\hat{Var}(v_1)]_0}$	$(v_2)_i / \sqrt{[\hat{Var}(v_2)]_0}$
18	<b>-1,4524</b>	-0,0809
19	0,3311	-0,4250
19	<b>-1,0514</b>	1,3836
20	0,5530	-0,5654
20	0,1288	0,4152
21	<b>-1,1524</b>	1,1128
21	<b>-1,1412</b>	0,0126
22	<b>-1,4256</b>	1,4346
22	<b>-0,6574</b>	-0,7431
23	<b>-1,3583</b>	1,4030
23	<b>-0,1483</b>	-1,1888
24	<b>-0,7258</b>	0,7348
24	<b>-0,2723</b>	-0,4410
25	<b>1,9768</b>	<b>-1,9814</b>
25	1,0032	0,9382
26	0,4584	-0,4183
26	0,7779	-0,3309
27	0,6324	-0,6015
27	0,7532	-0,1346
28	<b>-0,1868</b>	0,1920
28	<b>-0,0345</b>	-0,1493
29	<b>-0,2397</b>	0,2690
29	0,2549	-0,4925
30	<b>-1,2062</b>	1,2371
30	-0,2399	-0,9468
31	<b>-1,6281</b>	1,5954
31	<b>-1,3074</b>	-0,2888
32	1,0383	-1,0361
32	0,5853	0,4341
33	<b>-0,1648</b>	0,2045
33	0,4332	-0,5981
34	<b>-0,9889</b>	0,9534
34	<b>-0,9984</b>	0,0301
35	0,9476	-0,9653
35	0,2739	0,6579

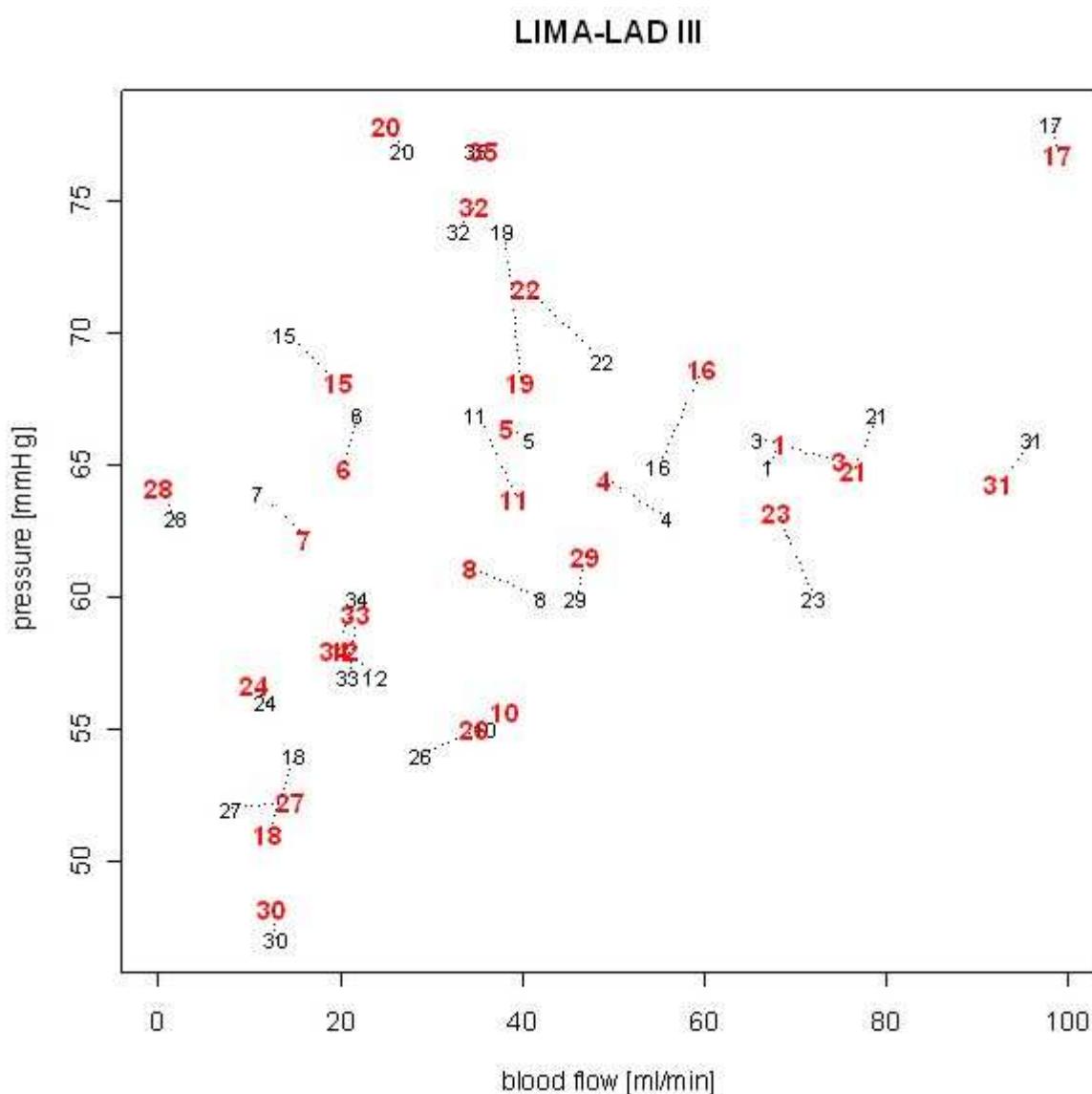
Po vyloučení našich „podezřelých bodů“ z výpočtu parametrů v modelu „Křišťálová koule“ získáme výsledky uvedené v tab. 7.

Srovnáme-li hodnoty vypočtených parametrů v případě kompletní množiny dat a množiny dat s vyloučenými podezřelými body, je vidět, že nyní průtok po povolení svorky více závisí na průtoku před povolením svorky a mnohem méně na tlaku před povolením svorky. V případě tlaku po povolení svorky lze rovněž konstatovat, že se zvýraznil rozdíl v závislosti na průtoku (nižší než pro kompletní data) a tlaku krve (vyšší než pro kompletní data) před povolením svorky na aortě.

Předpověď průtoku krve bypassem pomocí statistických metod

Tab. 7. Odhadnuté parametry po vyloučení podezřelých „outlierů“ z datového souboru.

$\hat{\beta}_1$	-45,67
$\hat{\beta}_2$	-19,17
$\hat{\beta}_3$	1,04
$\hat{\beta}_4$	0,50
$\hat{\beta}_5$	0,05
$\hat{\beta}_6$	1,24



Obr. 3. Předpovězené a naměřené hodnoty tlaku krve a průtoku krve bypassem po povolení svorky na aortě po vyloučení „podezřelých bodů“.

Předpověď průtoku krve bypassem pomocí statistických metod

Již na obr. 3, kde jsou tučně červeně vyznačeny odhady průtoku a tlaku po povolení svorky a černě odpovídající naměřené hodnoty, je vidět, že se po vyloučení „podezřelých bodů“ zmenšila rezidua. Potvrzuje to i tabulka reziduí.

Tab. 8. Tabulka reziduí po vyloučení podezřelých „outlierů“ z datového souboru.

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
1	-1,40	-0,86
3	-8,74	0,83
4	6,91	-1,50
5	2,66	-0,45
6	1,55	2,04
7	-5,27	1,73
8	7,68	-1,15
10	-2,13	-0,67
11	-4,31	3,19
12	3,48	-0,98

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
15	-5,96	1,80
16	-4,77	-3,69
17	-0,88	1,15
18	2,84	2,99
19	-1,79	5,79
20	2,00	-0,97
21	2,60	2,13
22	8,55	-2,76
23	4,09	-3,29
24	1,35	-0,64

i	$\Delta v_{i,1}$	$\Delta v_{i,2}$
26	-5,90	-0,99
27	-6,69	-0,22
28	1,94	-1,19
29	-0,87	-1,60
30	0,54	-1,13
31	3,57	1,65
32	-1,90	-0,86
33	-0,70	-2,34
34	2,42	1,95
35	-0,86	0,04

## 2.4 Nalezení „leverage“ bodů v modelu „Křišťálová koule“

V modelu „Křišťálová koule“ z kapitoly 1.2 – vztahy (2), (3) tohoto článku budeme hledat body, které nejvíce ovlivňují výsledné odhadnuté hodnoty tlaku a průtoku krve.

Hledáme maximum absolutních hodnot derivací výrazů (38) a (39), tedy maxima:

$$\left| \frac{\partial \mathbf{\Pi}_1 \mathbf{Z}_i}{\partial \mathbf{Z}_j} \right| \quad \left| \frac{\partial \mathbf{\Pi}_1}{\partial \mathbf{Z}_{i,j}} \right| \quad \text{a} \quad \left| \frac{\partial \hat{\mathbf{Z}}_i}{\partial \mathbf{Z}_j} \right| \quad \left| \frac{\partial \mathbf{\Pi}_2}{\partial \mathbf{Z}_{i,j}} \right|. \quad (7)$$

Pro naše data bez vyloučení podezřelých „outlierů“ získáme výpočtem jako „leverage“ pro odhady přírušků  $\delta\mu$  a  $\delta\nu$  i parametrů  $\beta$  bod s indexem 2.

Po vyloučení podezřelých „outlierů“ obdržíme jako „leverage“ pro odhady přírušků  $\delta\mu$  a  $\delta\nu$  bod s indexem 17 a pro odhad přírušků parametrů  $\beta$  bod s indexem 30.

## Závěr

Použijeme-li pro odhad tlaku a průtoku krve po povolení svorky během CABG operace klasický lineární model získáme méně přesné odhady než v případě modelu „Křišťálové koule“. Když navíc vyloučíme podezřelé „outlier“y, výsledné odhady se ještě více zpřesní. To nás vede k hypotéze, že pacienti tvoří několik skupin, přičemž pomocí navrhovaného algoritmu vypočteme určité hodnoty parametrů pro jednu skupinu pacientů a obecně odlišné hodnoty parametrů pro jinou skupinu. Nesprávné zařazení pacientů do jiné skupiny by se tak projevilo ve formě „outlier“u. Aby se tato hypotéza potvrdila, je zapotřebí provést hlubší analýzu většího souboru pacientů a zahrnout mezi zkoumané parametry ty, které by mohly

### Předpověď průtoku krve bypassem pomocí statistických metod

průtok krve bypassem před povolením svorky a po povolení svorky na aortě ovlivňovat a tak vytvářet výše zmíněné skupiny. Nabízí se zde např. procento stenózy v cílovém koronárním řečišti, ejekční frakce, FFT poměr, použitý štěp (LIMA, RIMA, SVG), doba dilatace, štěp jako pedikl/skelet apod. [2], [3], [12].

V této orientační etapě výzkumu byla použitá velmi zjednodušená struktura kovarianční matici. V dalších etapách je nutno prozkoumat kovarianční matici respektující různost disperzí u měření průtoku krve a u měření tlaku krve.

Numerické výsledky získané navrhovanou metodou „Křišťálová koule“ naznačují reálnou možnost předpovídat hodnoty průtoku a tlaku krve po povolení svorky na aortě. Před uvedením této metody do praxe je potřebné ještě provést řadu výpočtů na rozsáhlejších souborech dat.

## Apendix: Model „Křišťálová koule“

Model vztahu mezi průtokem a tlakem v době naložené svorky na aortě (LIMA-LAD I) a po povolení svorky (LIMA-LAD III) předpokládáme ve tvaru (značení, viz kapitola 1.1):

$$E(\zeta) = [\square_{1,1}, \square_{1,2}, \dots, \square_{n,1}, \square_{n,2}, \square_{1,1}, \square_{1,2}, \dots, \square_{n,1}, \square_{n,2}], \quad \text{Var}(\zeta) = \Sigma^2 \begin{pmatrix} 2n, \Theta & \mathbf{I} \\ \mathbf{I} & 2n, 2n \end{pmatrix}, \quad (8)$$

přičemž  $E(\zeta)$  splňuje podmínky:

$$\begin{pmatrix} E(\square_i) \\ E(\square_i) \end{pmatrix} = \begin{pmatrix} \square_1 \\ \square_2 \end{pmatrix} \begin{pmatrix} \square_3 & \square_4 \\ \square_5 & \square_6 \end{pmatrix} \begin{pmatrix} E(x_i) \\ E(y_i) \end{pmatrix}, \quad i=1,2,\dots,n. \quad (9)$$

Poznámka:  $E(\zeta)$  a  $\text{Var}(\zeta)$  značí střední hodnotu a varianci náhodného vektoru  $\zeta$ .

### A.1 Odhad parametrů modelu

Úlohou je na základě naměřených dat odhadnout hodnoty parametrů modelu tzn.  $\mu_{1,1}, \mu_{1,2}, \dots, \mu_{n,1}, \mu_{n,2}, \nu_{1,1}, \nu_{1,2}, \dots, \nu_{n,1}, \nu_{n,2}$  a zejména  $\beta_1, \dots, \beta_6$ , přičemž pro odhadnutelnost parametrů musí být počet pozorování plus počet podmínek vždy větší než počet parametrů, tj. v našem případě  $4n+2n > 4n+6$ , protože měříme 4 parametry u každého pacienta. K řešení úkolu se použije linearizace modelu. Model je totiž nelineární, neboť se v něm objevují součiny parametrů  $\beta$  a  $\mu$ . Po linearizaci použijeme odhadovacích algoritmů z modelu přímého neúplného měření vektorového parametru se systémem podmínek ze čtvrté kapitoly knihy Statistika a metrologie [14] a z knihy Statistical models with Linear Structures [15].

Předpověď průtoku krve bypassem pomocí statistických metod

Po linearizaci získáme model nepřímého měření neúplného vektorového parametru ve tvaru:

$$\zeta \begin{pmatrix} \square_{1,1}^{(0)} \\ \square_{1,2}^{(0)} \\ \vdots \\ \square_{n,1}^{(0)} \\ \square_{n,2}^{(0)} \end{pmatrix} \sim_{140} \begin{pmatrix} \square_{1,1}^{(0)} \\ \square_{1,2}^{(0)} \\ \vdots \\ \square_{n,1}^{(0)} \\ \square_{n,2}^{(0)} \end{pmatrix}; \square^2 \mathbf{I} \quad (10)$$

s podmínkami

$$\mathbf{b} \cdot \mathbf{B}_1 \begin{pmatrix} \square_{1,1}^{(0)} \\ \square_{1,2}^{(0)} \\ \vdots \\ \square_{n,1}^{(0)} \\ \square_{n,2}^{(0)} \end{pmatrix} \mathbf{B}_2 \begin{pmatrix} \square_1 \\ \vdots \\ \square_6 \end{pmatrix} = 0, \quad (11)$$

kde

$$\mathbf{b} = \begin{pmatrix} \square \square_{1,1}^{(0)} & \square_1^{(0)} & \square_3^{(0)} \square_{1,1}^{(0)} & \square_4^{(0)} \square_{1,2}^{(0)} \\ & \vdots & & \\ \square \square_{i,1}^{(0)} & \square_1^{(0)} & \square_3^{(0)} \square_{i,1}^{(0)} & \square_4^{(0)} \square_{i,2}^{(0)} \\ \square \square_{i,2}^{(0)} & \square_2^{(0)} & \square_4^{(0)} \square_{i,1}^{(0)} & \square_6^{(0)} \square_{i,2}^{(0)} \\ & \vdots & & \\ \square \square_{i,2}^{(0)} & \square_2^{(0)} & \square_4^{(0)} \square_{i,1}^{(0)} & \square_6^{(0)} \square_{i,2}^{(0)} \end{pmatrix}, \quad (12)$$

Matice  $\mathbf{B}_1$  je typu  $2n \times 4n$  a má následující tvar (svislá čára odděluje prvních  $2n$  sloupců):

$$\left( \begin{array}{cccccc|cccccc|cccccc} \square_3^{(0)} & \square_4^{(0)} & 0 & 0 & 0 & \cdots & 0 & 0 & 0 & \square 1 & 0 & 0 & 0 & 0 & \cdots & 0 & 0 & 0 \\ \square_5^{(0)} & \square_6^{(0)} & 0 & 0 & 0 & \cdots & 0 & 0 & 0 & 0 & \square 1 & 0 & 0 & 0 & \cdots & 0 & 0 & 0 \\ 0 & 0 & \square_3^{(0)} & \square_4^{(0)} & 0 & \cdots & 0 & 0 & 0 & 0 & 0 & \square 1 & 0 & 0 & \cdots & 0 & 0 & 0 \\ 0 & 0 & \square_5^{(0)} & \square_6^{(0)} & 0 & \cdots & 0 & 0 & 0 & 0 & 0 & 0 & \square 1 & 0 & \cdots & 0 & 0 & 0 \\ \cdots & \cdots \\ 0 & 0 & 0 & 0 & 0 & \cdots & 0 & \square_3^{(0)} & \square_4^{(0)} & 0 & 0 & 0 & 0 & 0 & \cdots & 0 & \square 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & \cdots & 0 & \square_5^{(0)} & \square_6^{(0)} & 0 & 0 & 0 & 0 & 0 & \cdots & 0 & 0 & \square 1 \end{array} \right). \quad (13)$$

Předpověď průtoku krve bypassem pomocí statistických metod

Matice  $\mathbf{B}_2$  je typu  $2n \times 6$  a má tvar:

$$\begin{pmatrix} 1 & 0 & \square_{1,1}^{(0)} & \square_{1,2}^{(0)} & 0 & 0 \\ 0 & 1 & 0 & 0 & \square_{1,1}^{(0)} & \square_{1,2}^{(0)} \\ 1 & 0 & \square_{2,1}^{(0)} & \square_{2,2}^{(0)} & 0 & 0 \\ 0 & 1 & 0 & 0 & \square_{2,1}^{(0)} & \square_{2,2}^{(0)} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & \square_{n,1}^{(0)} & \square_{n,2}^{(0)} & 0 & 0 \\ 0 & 1 & 0 & 0 & \square_{n,1}^{(0)} & \square_{n,2}^{(0)} \end{pmatrix}. \quad (14)$$

Odhad parametrů  $\delta\mu_{i,j}$ ,  $\delta\nu_{i,j}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  a  $\delta\beta_i$ ,  $i=1,2,\dots,6$  získáme minimalizací funkce:

$$\boxed{\square \square \square \square \square \square} \left( \zeta \boxed{\begin{pmatrix} \square_{1,1}^{(0)} \\ \vdots \\ \square_{n,2}^{(0)} \end{pmatrix}} \square \boxed{\begin{pmatrix} \square \square \square_{1,1} \\ \vdots \\ \square \square \square_{n,2} \end{pmatrix}} \right) \left( \zeta \boxed{\begin{pmatrix} \square_{1,1}^{(0)} \\ \vdots \\ \square_{n,2}^{(0)} \end{pmatrix}} \square \boxed{\begin{pmatrix} \square \square \square_{1,1} \\ \vdots \\ \square \square \square_{n,2} \end{pmatrix}} \right) \quad (15)$$

při splnění podmínek (11).

Řešení, po zdlouhavém výpočtu, dostáváme ve tvaru:

$$\begin{pmatrix} \square \square \square_{1,1} \\ \square \square \square_{1,2} \\ \vdots \\ \square \square \square_{n,1} \\ \square \square \square_{n,2} \\ \square \square \square_{1,1} \\ \square \square \square_{1,2} \\ \vdots \\ \square \square \square_{n,1} \\ \square \square \square_{n,2} \end{pmatrix} = \mathbf{Z} \square \mathbf{B}_1 \square \mathbf{B}_1 \mathbf{B}_1 \square \mathbf{B}_2 \mathbf{B}_2 \square^{-1} \square \mathbf{B}_1 \mathbf{B}_1 \square \mathbf{B}_2 \mathbf{B}_2 \square^{-1} \mathbf{B}_2 \square \mathbf{B}_1 \mathbf{B}_1 \square \mathbf{B}_2 \mathbf{B}_2 \square^{-1} \mathbf{B}_2 \square \mathbf{B}_1 \mathbf{B}_1 \square \mathbf{B}_2 \mathbf{B}_2 \square^{-1} \mathbf{B}_1 \mathbf{Z}. \quad (16)$$

$$\boxed{\hat{\square \square \square_1}, \hat{\square \square \square_2}, \hat{\square \square \square_3}, \hat{\square \square \square_4}, \hat{\square \square \square_5}, \hat{\square \square \square_6}} \quad \boxed{\mathbf{B}_2 \square \mathbf{B}_1 \mathbf{B}_1 \square \mathbf{B}_2 \mathbf{B}_2 \square^{-1} \mathbf{B}_2 \square \mathbf{B}_1 \mathbf{B}_1 \square \mathbf{B}_2 \mathbf{B}_2 \square^{-1} \mathbf{B}_1 \mathbf{Z}} \quad (17)$$

kde

$$\mathbf{Z} \square \boxed{x_1 \square \square_{1,1}^{(0)}, y_1 \square \square_{1,2}^{(0)}, \dots, x_n \square \square_{n,1}^{(0)}, y_n \square \square_{n,2}^{(0)}, \square_1 \square \square_{1,1}^{(0)}, \square_1 \square \square_{1,2}^{(0)}, \dots, \square_n \square \square_{n,1}^{(0)}, \square_n \square \square_{n,2}^{(0)}}. \quad (18)$$

Předpověď průtoku krve bypassem pomocí statistických metod

Toto řešení považujeme za výsledek prvního iteračního kroku, tzn.:

$$\begin{pmatrix} \hat{\Delta}_{1,1}^{(1)} \\ \vdots \\ \hat{\Delta}_{n,2}^{(1)} \end{pmatrix} \begin{pmatrix} \hat{\square}_{1,1}^{(0)} \\ \vdots \\ \hat{\square}_{n,2}^{(0)} \end{pmatrix} \begin{pmatrix} \hat{\square\Delta}_{1,1}^{(1)} \\ \vdots \\ \hat{\square\Delta}_{n,2}^{(1)} \end{pmatrix}; \begin{pmatrix} \hat{\square}_1^{(1)} \\ \vdots \\ \hat{\square}_6^{(1)} \end{pmatrix} \begin{pmatrix} \hat{\square}_1^{(0)} \\ \vdots \\ \hat{\square}_6^{(0)} \end{pmatrix} \begin{pmatrix} \hat{\square\hat{\square}}_1^{(1)} \\ \vdots \\ \hat{\square\hat{\square}}_6^{(1)} \end{pmatrix} \quad (19)$$

a v dalším iteračním kroku použijeme namísto vektorů

$$\begin{pmatrix} \mu_{1,1}^{(0)} \\ \vdots \\ \nu_{n,2}^{(0)} \end{pmatrix}, \begin{pmatrix} \beta_1^{(0)} \\ \vdots \\ \beta_6^{(0)} \end{pmatrix} \text{ vektory } \begin{pmatrix} \hat{\mu}_{1,1}^{(1)} \\ \vdots \\ \hat{\nu}_{n,2,korig}^{(1)} \end{pmatrix}, \begin{pmatrix} \hat{\beta}_1^{(1)} \\ \vdots \\ \hat{\beta}_6^{(1)} \end{pmatrix}.$$

Výsledné odhady lze tedy zapsat ve tvaru (kde  $k$  značí  $k$ -tou iteraci):

$$\begin{array}{ll} \hat{\Delta}_{i,1}^{(k)} & \hat{\Delta}_{i,1}^{(k \square 1)} & \hat{\square\Delta}_{i,1}^{(k)} \\ \hat{\Delta}_{i,2}^{(k)} & \hat{\Delta}_{i,2}^{(k \square 1)} & \hat{\square\Delta}_{i,2}^{(k)} \end{array} \quad \mathbf{a} \quad \begin{array}{ll} \hat{\square}_{i,1}^{(k)} & \hat{\square}_{i,1}^{(k \square 1)} & \hat{\square\hat{\square}}_{i,1}^{(k)} \\ \hat{\square}_{i,2}^{(k)} & \hat{\square}_{i,2}^{(k \square 1)} & \hat{\square\Delta}_{i,2}^{(k)} \end{array} \quad i=1,2,\dots,n,$$

$$\hat{\square}_i^{(k)} \quad \hat{\square}_i^{(k \square 1)} \quad \hat{\square\hat{\square}}_i^{(k)} \quad i=1,2,\dots,6,$$
(20)

kde  $\hat{\Delta}_{i,j}^{(k \square 1)}$ ,  $\hat{\Delta}_{i,j}^{(k \square 1)}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  a  $\hat{\square}_{i,j}^{(k \square 1)}$ ,  $i=1,2,\dots,6$  jsou odhadы parametrů  $\nu$  a  $\beta$  z předchozí iterace.

Vektor  $\mathbf{Z}^{(k)}$  v  $k$ -tému iteračnímu kroku vypadá takto:

$$[\mathbf{Z}^{(k)}]^T = [x_1 \square \hat{\Delta}_{1,1}^{(k \square 1)}, y_1 \square \hat{\Delta}_{1,2}^{(k \square 1)}, \dots, x_n \square \hat{\Delta}_{n,1}^{(k \square 1)}, y_n \square \hat{\Delta}_{n,2}^{(k \square 1)}, \hat{\square}_1 \square \hat{\Delta}_{1,1}^{(k \square 1)}, \hat{\square}_1 \square \hat{\Delta}_{1,2}^{(k \square 1)}, \dots, \hat{\square}_n \square \hat{\Delta}_{n,1}^{(k \square 1)}, \hat{\square}_n \square \hat{\Delta}_{n,2}^{(k \square 1)}] \quad (21)$$

Při iteračním výpočtu se odhady parametrů  $\nu$  korigují pomocí vztahů, které vycházejí ze snahy minimalizovat vliv nelinearity:

$$\begin{array}{ll} \hat{\Delta}_{i,1,korig.}^{(k)} & \hat{\Delta}_{i,2,korig.}^{(k)} \\ \hat{\Delta}_{i,1,korig.}^{(k)} & \hat{\Delta}_{i,2,korig.}^{(k)} \end{array} \quad \begin{array}{ll} \hat{\square}_1^{(k)} & \hat{\square}_2^{(k)} \\ \hat{\square}_2^{(k)} & \hat{\square}_5^{(k)} \end{array} \quad \begin{array}{ll} \hat{\square\Delta}_{3,1}^{(k)} & \hat{\square\Delta}_{4,1}^{(k)} \\ \hat{\square\Delta}_{5,i,1}^{(k)} & \hat{\square\Delta}_{6,i,1}^{(k)} \end{array} \quad \begin{array}{ll} \hat{\square\Delta}_{4,2}^{(k)} & \hat{\square\Delta}_{5,2}^{(k)} \\ \hat{\square\Delta}_{6,2}^{(k)} & \end{array} \quad i=1,2,\dots,n, k=0,1,2\dots$$
(22)

Pro nultou (počáteční) iteraci se volí

$$\hat{\square}_{i,1}^{(0)} = x_i, \quad \hat{\square}_{i,2}^{(0)} = y_i \quad \text{pro } i=1,2,\dots,n \quad (23)$$

Předpověď průtoku krve bypassem pomocí statistických metod

a parametry  $v_{i,j}^{(0)}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  a  $\beta_i^{(0)}$ ,  $i=1,2,\dots,6$  se vypočítají z rovnic sestavených pro tři body ležící v rovině  $xy$  na okraji pole, tj. pro jisté souřadnice:

$$\begin{pmatrix} x_{i1} \\ y_{i1} \end{pmatrix}, \begin{pmatrix} x_{i2} \\ y_{i2} \end{pmatrix}, \begin{pmatrix} x_{i3} \\ y_{i3} \end{pmatrix} \text{ a jim odpovídající body } \begin{pmatrix} \xi_{i1} \\ \eta_{i1} \end{pmatrix}, \begin{pmatrix} \xi_{i2} \\ \eta_{i2} \end{pmatrix}, \begin{pmatrix} \xi_{i3} \\ \eta_{i3} \end{pmatrix}.$$

Pro tyto body budeme požadovat, aby byla splněna soustava rovnic:

$$\begin{pmatrix} 1 & 0 & x_{i1} & y_{i1} & 0 & 0 \\ 0 & 1 & 0 & 0 & x_{i1} & y_{i1} \\ 1 & 0 & x_{i2} & y_{i2} & 0 & 0 \\ 0 & 1 & 0 & 0 & x_{i2} & y_{i2} \\ 1 & 0 & x_{i3} & y_{i3} & 0 & 0 \\ 0 & 1 & 0 & 0 & x_{i3} & y_{i3} \end{pmatrix} \begin{pmatrix} \square_1^{(0)} \\ \square_2^{(0)} \\ \square_3^{(0)} \\ \square_4^{(0)} \\ \square_5^{(0)} \\ \square_6^{(0)} \end{pmatrix} = \begin{pmatrix} \square_{i1} \\ \square_{i2} \\ \square_{i3} \\ \square_{i2} \\ \square_{i3} \\ \square_{i3} \end{pmatrix} \quad (24)$$

Řešením této soustavy jsou parametry  $\beta_i^{(0)}$ ,  $i=1,2,\dots,6$ , s jejichž pomocí se vypočítají nulté iterace parametrů  $v_{i,j}^{(0)}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  na základě vztahů:

$$\begin{array}{cccccc} \square_{i,1}^{(0)} & \square_1^{(0)} & \square_3^{(0)} & \square_{i,1}^{(0)} & \square_4^{(0)} & \square_{i,2}^{(0)} \\ \square_{i,2}^{(0)} & \square_2^{(0)} & \square_5^{(0)} & \square_{i,1}^{(0)} & \square_6^{(0)} & \square_{i,2}^{(0)} \end{array} \quad i=1,2,\dots,n, \quad (25)$$

přičemž parametry  $\mu_{i,j}^{(0)}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  se volí takto:

$$\square_{i,1}^{(0)} \quad x_i \quad \text{a} \quad \square_{i,2}^{(0)} \quad y_i. \quad (26)$$

Iterační výpočet zastavíme, když lze odhadu přírůstků  $\delta\mu_{i,j}$ ,  $\delta v_{i,j}$ ,  $i=1,2,\dots,n$ ,  $j=1,2$  a  $\delta\beta_i$ ,  $i=1,2,\dots,6$  zanedbat, tj. platí podmínka pro ukončení iteračního výpočtu:

$$|\hat{\square}_{i,j} - \hat{\square}_{i,j,korig.}| \leq \varepsilon \quad i=1,2,\dots,n, j=1,2, \quad (27)$$

kde  $\varepsilon$  je předem daná konstanta.

Volbou  $i1=14$ ,  $i2=17$ ,  $i3=18$  pro výpočet nulté iterace a konstanty  $\varepsilon=0,1$  získáme pro naše data po šesti iteračních krocích pro hledané parametry  $\beta_i$  numerické výsledky uvedené v tabulce 4.

Kovarianční matice pro tyto parametry se vypočte dle následujícího vzorce:

$$\hat{V}ar \begin{pmatrix} \hat{\square}_1 \\ \hat{\square}_2 \\ \hat{\square}_3 \\ \hat{\square}_4 \\ \hat{\square}_5 \\ \hat{\square}_6 \end{pmatrix} = \hat{\sigma}^2 \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{I} \quad (28)$$

kde

$$\hat{\sigma}^2 = \frac{\mathbf{v}_1^\top \mathbf{v}_1 + \mathbf{v}_2^\top \mathbf{v}_2}{64} \quad (29)$$

přičemž

$64 =$  počet měření (140) + počet podmínek (70) - počet parametrů (140+6)

a

$$\begin{pmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{pmatrix} = \begin{pmatrix} x_1 & \hat{\square}_{1,1} \\ \vdots & \vdots \\ y_n & \hat{\square}_{n,2} \\ \hat{\square}_1 & \hat{\square}_{1,1} \\ \vdots & \vdots \\ \hat{\square}_n & \hat{\square}_{n,2} \end{pmatrix} \quad (30)$$

Pro naše data vypadá numerický výsledek takto:

$$\hat{\sigma}^2 = 46.52034,$$

$$\hat{V}ar \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \\ \hat{\beta}_4 \\ \hat{\beta}_5 \\ \hat{\beta}_6 \end{pmatrix} = \begin{pmatrix} 539.914 & 225.030 & -0.840 & -7.528 & -0.350 & -3.137 \\ 225.030 & 326.432 & -0.350 & -3.137 & -0.508 & -4.551 \\ -0.840 & -0.350 & 0.008 & 0.005 & 0.003 & 0.002 \\ -7.528 & -3.137 & 0.005 & 0.110 & 0.002 & 0.046 \\ -0.350 & -0.508 & 0.003 & 0.002 & 0.005 & 0.003 \\ -3.137 & -4.551 & 0.002 & 0.046 & 0.003 & 0.066 \end{pmatrix}$$

Předpověď průtoku krve bypassem pomocí statistických metod

## A.2 Nalezení odlehlých („outlier“) bodů v modelu „Křišťálová koule“

Jakmile je iterační proces zastaven (v  $k$ -tému kroku), potom pro náš vektor reziduú  $\mathbf{Z}$  ze vztahu (21) určitě platí:

$$\mathbf{Z}^{(k)} \left( \begin{matrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{matrix} \right) \left( \begin{matrix} x_1 & \square \hat{\Delta}_{1,1}^{(k)} \\ \vdots & \vdots \\ y_n & \square \hat{\Delta}_{n,2}^{(k)} \\ \square_1 & \square \hat{\Delta}_{1,1}^{(k)} \\ \vdots & \vdots \\ \square_n & \square \hat{\Delta}_{n,2}^{(k)} \end{matrix} \right) \left( \begin{matrix} x_1 \\ \vdots \\ y_n \\ \square_1 \\ \vdots \\ \square_n \end{matrix} \right) \square \left( \begin{matrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{matrix} \right) \left( \begin{matrix} \hat{\boldsymbol{\mu}}^{(k)} \\ \hat{\mathbf{v}}^{(k)} \end{matrix} \right) \left( \begin{matrix} x_1 \\ \vdots \\ y_n \\ \square_1 \\ \vdots \\ \square_n \end{matrix} \right) \square \left( \begin{matrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{matrix} \right) \left( \begin{matrix} \boldsymbol{\mu}^{(k-1)} & \delta\hat{\boldsymbol{\mu}}^{(k)} \\ \mathbf{v}^{(k-1)} & \delta\hat{\mathbf{v}}^{(k)} \end{matrix} \right)$$

$$\left( \begin{matrix} x_1 \\ \vdots \\ y_n \\ \square_1 \\ \vdots \\ \square_n \end{matrix} \right) \square \boldsymbol{\mu}^{(k-1)}$$

$$\left( \begin{matrix} x_1 \\ \vdots \\ y_n \\ \square_1 \\ \vdots \\ \square_n \end{matrix} \right) \square \mathbf{v}^{(k-1)}$$
(31)

$$\square \left[ \mathbf{Z}^{(k-1)} \square \mathbf{B}_1^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{Z}^{(k-1)} \right]$$

Dosadíme-li místo prvního členu výrazu na posledním řádku vektor  $\mathbf{Z}$  z poslední iterace, pak se celý výraz dá upravit na výsledný tvar:

$$\left( \begin{matrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{matrix} \right) \square \mathbf{B}_1^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \square \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{Z}^{(k-1)}$$
(32)

Potom, označíme-li

$$\mathbf{T} = \left[ \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \right] \quad (33)$$

můžeme psát:

$$Var \left( \begin{matrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{matrix} \right) = \square^2 \mathbf{B}_1^\top \mathbf{T} \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{T} \mathbf{B}_1 = \begin{pmatrix} Var(\mathbf{v}_1) & cov(\mathbf{v}_1, \mathbf{v}_2) \\ cov(\mathbf{v}_2, \mathbf{v}_1) & Var(\mathbf{v}_2) \end{pmatrix} = \square^2 \begin{pmatrix} \mathbf{Q}_{11} & \mathbf{Q}_{12} \\ \mathbf{Q}_{21} & \mathbf{Q}_{22} \end{pmatrix} \quad (34)$$

Pro každé  $i=1,2,\dots,2n$  vypočítáme výrazy:

$$|\mathbf{v}_1| / \sqrt{Var(\mathbf{v}_1)}_i \quad |\mathbf{v}_1| / \sqrt{\square^2 \mathbf{Q}_{11}}_i$$

$$\quad \quad \quad (35)$$

$$|\mathbf{v}_2| / \sqrt{Var(\mathbf{v}_2)}_i \quad |\mathbf{v}_2| / \sqrt{\square^2 \mathbf{Q}_{22}}_i$$

Předpověď průtoku krve bypassem pomocí statistických metod

### A.3 Nalezení „leverage“ bodů v modelu „Křišťálová koule“

Jakmile je iterační proces zastaven (v k-tém kroku), použijeme rovnice pro přírůstky

$$\begin{pmatrix} \hat{\square}_{1,1} \\ \hat{\square}_{1,2} \\ \vdots \\ \hat{\square}_{n,1} \\ \hat{\square}_{n,2} \\ \hat{\square}_{1,1} \\ \hat{\square}_{1,2} \\ \vdots \\ \hat{\square}_{n,1} \\ \hat{\square}_{n,2} \end{pmatrix} \mathbf{Z}^{(k)} = \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{Z}^{(k)} \quad (36)$$

$$[\hat{\square}_1, \hat{\square}_2, \hat{\square}_3, \hat{\square}_4, \hat{\square}_5, \hat{\square}_6]^\top = \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{Z}^{(k)} \quad (37)$$

které přepíšeme do tvaru:

$$\Pi_1 \mathbf{Z}^{(k)} = \begin{pmatrix} \hat{\square}_{1,1} \\ \hat{\square}_{1,2} \\ \vdots \\ \hat{\square}_{n,1} \\ \hat{\square}_{n,2} \\ \hat{\square}_{1,1} \\ \hat{\square}_{1,2} \\ \vdots \\ \hat{\square}_{n,1} \\ \hat{\square}_{n,2} \end{pmatrix} \quad (38)$$

$$\mathbf{I} = \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{Z}^{(k)}$$

$$\Pi_2 \mathbf{Z}^{(k)} = [\hat{\square}_1, \hat{\square}_2, \hat{\square}_3, \hat{\square}_4, \hat{\square}_5, \hat{\square}_6]^\top = \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{B}_1^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_2 \mathbf{B}_2^\top \mathbf{B}_1 \mathbf{Z}^{(k)} \quad (39)$$

## Literatura

- [1] ÚZIS: Aktuální informace č. 14 - Kardiochirurgické operace (k dispozici na [http://www.uzis.cz/cz/archiv04/14\\_04.pdf](http://www.uzis.cz/cz/archiv04/14_04.pdf))

Předpověď průtoku krve bypassem pomocí statistických metod

- [2] Takami, Y., Ina, H.: A simple Method to Determine Anastomotic Quality of Coronary Artery Bypass Grafting in the Operating Room. *Cardiovascular Surgery*, Vol. 9, No. 5, pp. 499-503
- [3] Takami, Y., Ina, H.: Relation of Intraoperative Flow Measurement With Postoperative Quantitative Angiographic Assessment of Coronary Artery Bypass Grafting. *Ann Thorac Surg* 2001, 72:1270-4
- [4] D'Anconna, G., Karamanoukian, H. L., Bergsland, J.: Is Intraoperative measurement of coronary blood flow a good predictor of graft patency? *European Journal of Cardio-thoracic Surgery* 20 (2001), pp. 1075-1076
- [5] Shin, H. et al.: Intraoperative Assessment of Coronary Artery Bypass Graft: Transit-Time Flowmetry Versus Angiography. *Ann Thorac Surg* 2001, 72:1562-5
- [6] Louagie, Y. et al.: Pulsed Doppler Intraoperative Flow Assessment and Midterm Coronary Graft Patency. *Ann Thorac Surg* 1998, 66:1282-8
- [7] Louagie, Y., Brockmann, C., Gurné, O., Jamart, J.: Intraoperative Flow Measurements: Predictive Value for Postoperative Angiographic Follow-Up. In.: *Intraoperative Graft Patency Verification in Cardiac and Vascular Surgery* (ed. G. D'Ancona). Futura Publishing, Armonk (NY) 2001
- [8] Milnor, W. R.: *Hemodynamics*. Williams&Wilkins, Baltimore (USA) 1998
- [9] Hata, M. et al.: Midterm Results of Coronary Artery Bypass Graft Surgery With Internal Thoracic Artery Under Low Free-Flow Conditions. *Ann Thorac Surg* 2004, 78:477-80
- [10] Pagni, S. et al.: Factors Affecting Internal Mammary Artery Graft Survival: How Is Competitive Flow from a Patent Native Coronary Vessel a Risk Factor? *Journal of Surgical Research* 71 (1997), pp. 172-178
- [11] Sabik, J. F. et al.: Does Competitive Flow Reduce Internal Thoracic Artery Graft Patency? *Ann Thorac Surg* 2003, 76:1490-7
- [12] Pagni, S., Storey, J. et al.: ITA versus SVG: a comparison of instantaneous pressure and flow dynamics during competitive flow. *European Journal of Cardio-thoracic Surgery* 11 (1997), pp. 1086-1092
- [13] Lausten, J.: Transit-Time Flow Measurement: Principles and Clinical Applications. In.: *Intraoperative Graft Patency Verification in Cardiac and Vascular Surgery* (ed. G. D'Ancona). Futura Publishing, Armonk (NY) 2001
- [14] Kubáček, L., Kubáčková, L.: *Statistika a metrologie*. VUP Olomouc 2000
- [15] Kubáček, L., Kubáčková, L., Volaufová, J.: *Statistical Models with Linear Structures*. Veda, Bratislava 1995.