

Integrating Kernel Methods Into a Knowledge-based Approach to Evidence-based Medicine

Katharina Morik

CS Department, AI Unit
University of Dortmund
D-44221 Dortmund
Germany

Thorsten Joachims

CS Department, AI Unit
University of Dortmund
D-44221 Dortmund
Germany

Michael Imhoff

Surgical Department
Community Hospital Dortmund
D-44137 Dortmund
Germany

Peter Brockhausen

CS Department, AI Unit
University of Dortmund
D-44221 Dortmund
Germany

Stefan Rüping

CS Department, AI Unit
University of Dortmund
D-44221 Dortmund
Germany

Operational protocols are a valuable means for quality control. However, developing operational protocols is a highly complex and costly task. We present an integrated approach involving both intelligent data analysis and knowledge acquisition from experts that supports the development and validation of operational protocols. The aim is to lower development cost through the use of machine learning and at the same time ensure high quality standards for the protocol through empirical validation. We demonstrate our approach of integrating expert knowledge with data driven techniques based on our effort to develop an operational protocol for the hemodynamic system.

1 Introduction

An abundance of information is generated during the process of critical care. Much of this information can now be captured and stored using clinical information systems (CIS) that have become commercially available for use in intensive care over the last years. These systems provide for a complete medical documentation at the bedside and their clinical usefulness and efficiency has been shown repeatedly [6, 7, 11]. While databases with more than 2,000 separate patient-related variables are now available for further analysis [8], the multitude of variables presented at the bedside even without a CIS precludes medical judgement by humans. A physician may be confronted with more than 200 variables in the critically ill during a typical morning round [21]. We know, however, that even an experienced physician is often not able to develop a systematic response to any problem involving more than seven variables [18]. Moreover, humans are limited in their ability to estimate the degree of relatedness between only two variables [12]. This problem is most pronounced in the evaluation of the measurable effect of a therapeutic intervention. Personal bias, experience, and a certain expectation toward the respective intervention may distort an objective judgement [4]. These arguments motivate the use of decision support systems.

Clinical decision support aims at providing health care professionals with therapy guidelines directly at the bed-side. This should enhance the quality of clinical care, since the guidelines sort out high value practices from those that have little or no value. The goal of decision support is to supply the best recommendation under all circumstances [22]. The computerized protocol of care can take into account more aspects of the patient than a physician can accommodate. It is not disturbed by circumstances or hospital constraints. It bridges the gap between low-level numerical measurements (the level of the equipment) and high-level qualitative principles (the level of medical reasoning). While knowledge-based systems have mostly been applied for diagnosis and therapy planning (e.g. [25], [16]), some systems also aim at on-line patient monitoring [5, 17, 22]. Methods that have proved their value in handling low-frequency patient data are not applicable for on-line monitoring [17]. Quantitative measurements and qualitative

reasoning have to be integrated in a system that recommends interventions in real-time. The numerical measurements of the patients' vital signs have to be abstracted into qualitative terms of high abstraction. The aspect of time has to be handled both at the level of measurements and the level of expert knowledge [3, 14, 17, 25]. In the expert's reasoning, time becomes the relation between time intervals, abstracting from the exact duration of, e.g., an increasing heart rate, and focusing on tendencies of other parameters (e.g., cardiac output) within overlapping time intervals.

One of the big obstacles to the more frequent implementation of decision support systems is the tedious and time-consuming task of developing the knowledge base. The decision support system for respiratory care at the LDS Hospital, Salt Lake City, USA [22], for instance, has been developed in about 25 person years. The method of guideline development itself is not supported by a computer system. Mechanisms of temporal abstraction and reasoning presuppose manually designed models or ontologies [3, 17, 25]. Why not use techniques of knowledge discovery and statistical time series analysis in order to ease the process of guideline generation? Machine learning and statistical analysis have been applied in building-up diagnostical systems successfully (e.g., [15]).

We now want to exploit the huge amount of data for the development of guidelines for on-line monitoring. Our task is to build a decision support system for on-line hemodynamic monitoring in the critically ill. We do not aim at modeling the actual physician's behavior. Imitating the actual interventions made by physicians is not the goal. Actual behavior is influenced by the overall hospital situation, e.g., how long is the physician on duty, how many patients require attention at the same time. Machine learning from patients' data could lead to a knowledge base that mirrors such disturbing effects. Therefore, the learned decision rules have to be checked by additional rules about effects of drug and fluid administration. Our approach is to combine statistics, knowledge acquisition, and machine learning. Our aim is to develop a method for guideline generation that is faster and more reliable than current methods.

Data for statistical evaluation and learning can be provided by the CIS. However, the nature of the data is different from that gathered in controlled experiments. While a CIS in modern intensive care can take numerous measurements every minute, the values of some vital signs are sometimes recorded only once every hour. Other vital signs are recorded only for a subset of the patients. Hence, the overall high dimensional data space is sparsely populated. Moreover, the average time difference between intervention as charted and estimated hemodynamic effect can show a wide variation [10]. Even the automatic measurements can be noisy due to manipulation of measurement equipment, flushing of pressure transducers, or technical artifacts. In some cases, relevant demographic and diagnostic parameters may even not be recorded at all. In summary, we have a large amount of high dimensional, numerical time series data that contains missing values and noise. Using this data already at the stage of development of the decision support system stave off surprises at the stage of clinical experience as has been reported in [17, p. 572]: “The huge number of measurements classified as invalid is quite astonishing although it reflects the real clinical environments.”

In addition to problems of knowledge acquisition, we see a particular need for knowledge validation. It should be noted that many medical guidelines published today are neither evidence-based nor sufficiently validated against real patient data. The current procedure is to first develop the guideline, then represent it in a knowledge-based system, and finally to test it in clinical studies. In this “waterfall” process, unrealistic assumptions, mistakes, and flaws are recognized at a late stage. In contrast, our approach includes validation from the very beginning. Using a knowledge-based system early on supports the validation of the knowledge base at earlier stages. Inconsistencies within the knowledge base as well as a mismatch of rules and patient data are detected while developing the knowledge base. A mismatch may indicate that the model underlying the knowledge base is insufficient. Hence, applying the model to patient data helps to find errors in its design. A mismatch may also indicate a difference in the medical practices of the physician at the bed-side and the medical expert that helped to develop the knowledge base. Moreover, experts from different schools or countries can vary quite a bit in their behavior

and knowledge. Matching the formally modeled guidelines with patient data facilitates and focuses the knowledge-acquisition process.

In order to test our approach to using real clinical data for building and validating a knowledge base for on-line monitoring, we have constructed a system. Its overall architecture is shown in Figure 1. The patients' measurements are used to recommend an intervention and are abstracted with respect to their course over time. The recommendation of interventions constitutes a model of physician behavior. This asks for further validation. Therefore, a recommended intervention is checked by calculating its expected effects on the basis of medical knowledge. In this way, a qualitative assessment of a statistical prediction enhances the model of physician behavior in order to obtain a model of best practice. The medical knowledge constitutes a model of the patients' hemodynamic system. This model is validated with respect to past patients' data. In detail, the processes we have designed are:

Data abstraction: Given series of measurements of one vital sign of the patient, detect and possibly eliminate outliers and find level changes by good statistical practice. This abstracts the measurements to qualitative propositions with respect to a time interval, e.g. within time point 12 and time point 63, the heart rate remained about equal, from time point 63 to time point 69 it was increasing. We used the statistical time series techniques of ARMA modeling and phase space embedding [1,2,9]

Data-driven acquisition of state-action rules: Given the numerical data describing signs of the patient, his or her current medication, find the appropriate intervention. An intervention is formalized as increasing, decreasing or not changing the dose of a drug. The decision is made every minute. These rules were learned by the Support Vector Machine [26].

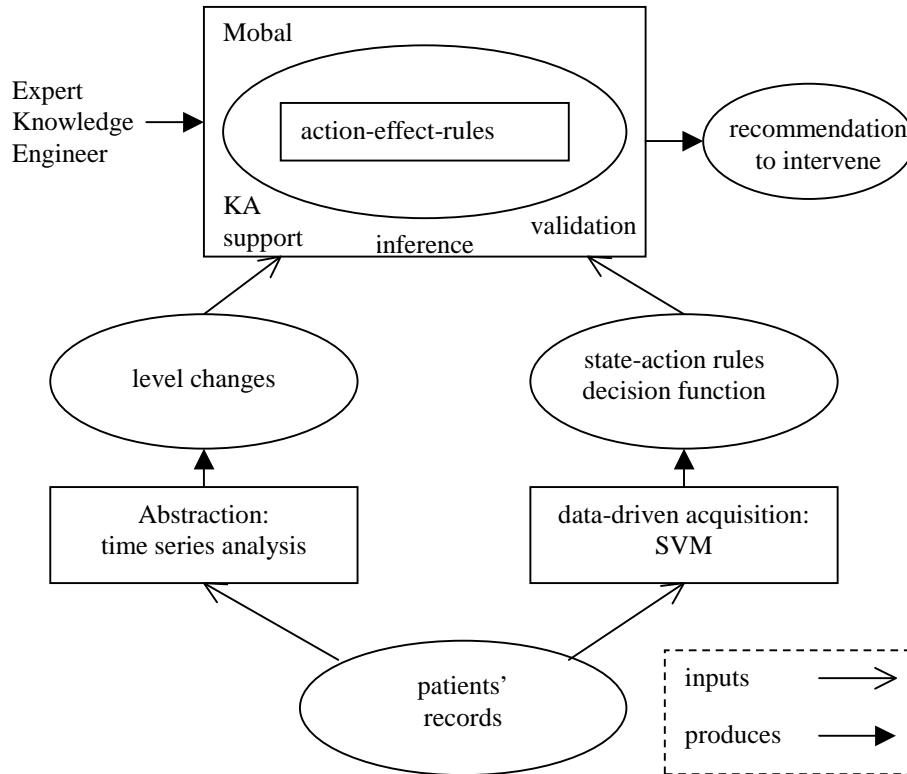


Figure 1. Overall system architecture.

Acquisition of medical knowledge: Given text book knowledge and explanations by an expert, represent the effects of substances in different dosages, relations between vital signs, and interrelations between different substances, and validate the knowledge on the basis of past patients' data. The knowledge acquisition and validation was supported by the MOBAL system [20].

validation of recommended interventions: Given

- the state of a patient described in qualitative terms,
- medical knowledge
- a sequence of interventions, and
- a current intervention,

find the effects of the current intervention on the patient. The derivation of effects is made for each intervention as forward inference within MOBAL. The effect should result in a stable state of the patient.

The outline of this chapter is as follows. Throughout the chapter we report on the continuous development of a decision support system for intensive care as performed at the city hospital and the university of Dortmund. We start with a description of the data acquisition process at the hospital and the resulting data set [11]. Section 3 shows, how we applied the support vector machine (SVM) to learn state-action rules. A short introduction to the MOBAL system [20] and its representation of medical knowledge leads to the issue of validation which is presented in section 4.

2 Data Acquisition and Data Set

2.1 Data Acquisition

Most variables are entered by hand at the bedside. For entities such as clinical observations, nursing procedures, therapeutic measures, medications, or orders it appears very unlikely that entry of these variables can be automated in the foreseeable future. Only 5-10% of all variables in a CIS are acquired automatically. This includes the majority of bedside devices, e.g. physiologic monitors, ventilators, infusion devices. Additional data is interfaced from the hospital information system (HIS), the laboratory (LIS) or the microbiology information systems, where the LIS represents the clinically most relevant set of data among these centralized information systems. Although device data account for a comparatively small number of variables, they can, depending on the sampling rate, generate large amounts of data.

Table 1. Overall attribute set for learning state-effect rules

16 demographic attributes	5 intensive care diagnoses	6 continuously infused drugs
11 vital signs	9 derived parameters	14 respiratory variables
37 intake/output variables	10 bolus drugs	10 laboratory tests

The data structure of a CIS shows a wide variety of different data types on different scales (nominal scales, e.g. sex, breathing sounds; ordinal scales, e.g. neurological scoring; absolute scales, e.g. vital signs), which are stored at different time intervals (ranging from seconds for vital signs to once during the length of stay for demographic data). Time intervals may also be regular or irregular.

For further analysis data must be structured, so that it can be subjected to statistical algorithms. Numeric data, e.g. vital signs, intake/output, is typically directly accessible for most applications. Free-text data, which traditionally makes up a large portion of medical documentation, cannot be statistically analyzed in any structured way. Therefore, free-text entries into a CIS should be avoided wherever possible. Qualitative information, such as clinical observations or interventions, should be documented in a strictly structured fashion with selection lists and menu items. This approach provides a consistent terminology throughout the entire medical institution. It is highly efficient and fast, especially for users not well trained in the use of computers and keyboards in particular. In clinical practice, with the stringent implementation of structured tabular documentation, it was possible to reduce the use of free-text notes by more than 90%. Structured qualitative data can, in contrast to free-text information, be directly exported for statistical analysis.

These general propositions also hold for the city hospital of Dortmund, a 1,900-bed tertiary referral center. There, all medication data of the 16-bed surgical intensive care unit was charted with a CIS, allowing the user one minute time resolution for all data. Moreover, data from bedside devices, e.g. patient monitors, is gathered automatically every minute.

Table 2. Best feature set for learning state-action rules using SVM.

Vital signs (measured every minute)	Continuously given drugs (changes charted at 1-min-resolution)	Demographic Attributes (charted once at admission)
Diastolic Arterial Pressure	Dobutamine	Broca-Index
Systolic Arterial Pressure	Adrenaline	Age
Mean Arterial Pressure	Glyceroltrinitrate	Body Surface Area
Heart Rate	Noradrenaline	Emergency Surgery y/n
Central Venous Pressure	Dopamine	
Diastolic Pulmonary Pressure	Nifedipine	
Systolic Pulmonary Pressure		
Mean Pulmonary Pressure		

2.2 Data Set

The entire database of intensive care patient records at the city hospital of Dortmund comprises about 2,000 different variables (attributes). Data from the CIS is selected through customizable data filters and copied into a standard relational database where it is accessible for further data analysis.

For this investigation, data was acquired from 148 consecutive critically ill patients (53 female, 95 male, mean age 64.1 years), who had pulmonary artery catheters for extended hemodynamic monitoring. Recording in one minute intervals, this amounts to 679,817 sets of observations.

From the original database 118 attributes in 9 groups were taken for learning state-action rules (Table 1).

Categorical attributes are broken down into a number of binary attributes, each taking the values $\{0,1\}$. Real valued parameters are either scaled so that all measurements lie in the interval $[0,1]$, or they are normalized by empirical mean and variance:

$$\text{norm}(X) = (X - \text{means}(X)) / \sqrt{\text{var}(X)} \quad (1)$$

We systematically evaluated a large number of plausible attribute sets using a train/test scheme on the learning task described in section 3.2. The set with the best performance is given in Table 2. These attributes

are actually the most important parameters of the patient according to expert judgement. Only the relevant attributes “Cardiac Output” and “Net Intake/Output” are missing, but they cannot be used as they are not continuously available.

We also experimented with different ways of incorporating the history of the patient. We tried:

- using only the last minute before the intervention
- using the last up to 10 minutes before the intervention
- using the averages of up to 60 minutes before the intervention
- combinations of these
- the state of the patient at the previous intervention

None of the more complex approaches gave significantly better results on the learning task in section 3.2 than just using the measurements from one minute before the intervention. All the feature selection experiments were done on the training set, leaving a separate test set to measure the results presented in this chapter.

Since each patient record covers several interventions, data from 148 patients gives us sufficiently large sets of examples. For learning state-action rules, we used a total of 1319 training and 473 test examples. For the rule validation we analyzed 8200 interventions corresponding to 27400 intervention-effect pairs.

2.3 Statistical Preprocessing

Given series of measurements of one vital sign of the patient, the goal of statistical data abstraction is to detect and possibly eliminate outliers and find level changes by good statistical practice. This abstracts the measurements to qualitative propositions with respect to a time interval, e.g., within time point 12 and time point 63, the heart rate remained about equal, from time point 63 to time point 69 it was increasing. We used an approach based on statistical time series analysis. Classical ARMA (autoregressive moving average) modeling [2] is applied with corresponding outlier- and level shift detection procedures using the new tool of a phase space embedding [1,9].

3 Data-driven Acquisition of State-Action Rules

3.1 Support Vector Machine

Support vector machines (SVMs) [26] represent a method to learn either binary classifiers or function approximators from examples. For a set of training examples they find the classification rule for which they can guarantee the lowest error rate on new observations. Each example consists of a vector (describing e.g. the state of a patient represented by the current measurements of blood pressures, heart rate, etc.) and its label (classification or functional value).

In their basic form, SVMs learn linear decision rules $h(\vec{o}) = \text{sign}(\vec{w} \cdot \vec{o} + b)$. The weight vector \vec{w} and the threshold b are the result of learning and describe a hyperplane. Observations are classified according to which side of the hyperplane they are located. A typical decision rule is given in Figure 2. During training, the SVM calculates the hyperplane so that it classifies most training examples correctly while keeping a large “margin” around the hyperplane. If the training data can be separated without error, the margin is the distance from the hyperplane to the closest training examples.

Since we will be dealing with very unbalanced numbers of positive and negative examples in the following, we introduce cost factors to be able to adjust the cost of false positives vs. false negatives. Training an SVM can now be translated into the following optimization problem:

$$\text{Minimize: } J(\vec{w}, b, \vec{\xi}) = \frac{1}{2} \vec{w} \cdot \vec{w} + C_+ \sum_{i: y_i=1} \xi_i + C_- \sum_{j: y_j=-1} \xi_j \quad (2)$$

$$\text{subject to: } \forall t \in \{1 \dots n\}: y_t [\vec{w} \cdot \vec{o}_t + b] \geq 1 - \xi_t \wedge \xi_t \geq 0 \quad (3)$$

Training error is represented by the variables ξ_i, ξ_j , while the margin is measured by $\vec{w} \cdot \vec{w}$. We solve this optimization problem in its dual formulation using SVMlight [13], extended to handle unsymmetrical cost-factors.

$$h_{\text{nitroun}}(\vec{o}) = \text{sign} \left(\begin{array}{l} 0.014 \\ 0.019 \\ -0.001 \\ -0.015 \\ -0.016 \\ 0.026 \\ 0.134 \\ -0.177 \\ -9.543 \\ -1.047 \\ -0.185 \\ 0.542 \\ -0.017 \\ 2.391 \\ 0.033 \\ 0.334 \\ 0.784 \\ 0.015 \end{array} \cdot \begin{array}{l} \textit{Artsys} \ 174.00 \\ \textit{Artdia} \ 86.00 \\ \textit{Artmn} \ 121.00 \\ \textit{CVP} \ 8.00 \\ \textit{HR} \ 79.00 \\ \textit{Papsys} \ 26.00 \\ \textit{Papdia} \ 13.00 \\ \textit{Papmn} \ 15.00 \\ \textit{Nifedipine} \ 0.00 \\ \textit{Noradrenaline} \ 0.00 \\ \textit{Dobutamine} \ 0.00 \\ \textit{Dopamine} \ 0.00 \\ \textit{Glyceroltrinitrate} \ 0.00 \\ \textit{Adrenaline} \ 0.00 \\ \textit{Age} \ 77.91 \\ \textit{Emerg} \ 0 \\ \textit{BSA} \ 1.79 \\ \textit{Broca} \ 1.02 \end{array} \right) - 4.368$$

Figure 2. Decision rule and an instantiation for predicting an intervention that increases the dosage of Glyceroltrinitrate.

3.2 Learning the Directions of Interventions

The first question we asked ourselves was: Given that we know the physician changed the dosage of some drug, can we learn when he increased the dosage and when he decreased the dosage based on the state of the patient? For each drug, examples are taken from the points in time where, in fact, the dosage changed. For all drugs, linear SVMs are trained on the problem “increase of dosage” ($y_t = 1$) vs. “decrease of dosage” ($y_t = -1$) using the attributes in Table 2 for describing the

Table 3. Accuracy in predicting the right direction of an intervention

Drug	Accuracy	StdErr
Dobutamine	83.6%	2.6%
Adrenaline	81.3%	3.7%
Glyceroltrinitrate	85.5%	3.0%
Noradrenaline	86.0%	5.2%
Dopamine	84.0%	7.3%
Nifedipine	86.8%	7.0%

state of the patient. The performance of the respective SVM on a previously untouched test set is given in Table 3.

To get an impression about how good these prediction accuracies are, we conducted an experiment with a physician. On a subset of 41 test examples we asked an expert to do the same task as the SVM for Dobutamine, given the same information about the state of the patient. In a blind test the physician predicted the same direction of dosage change as actually performed in 32 out of the 41 cases. On the same examples the SVM predicted the same direction of dosage change as actually performed in 34 cases, resulting in an essentially equivalent accuracy.

3.3 Learning When to Intervene

The previous experiment shows that SVMs can learn in how far drugs should be changed given the state the patient is in. In reality, the physician also has to decide when to intervene or just keep a dosage constant. This leads to the following three class learning problem. Given the state of the patient, should the dosage of a drug be increased, decreased or kept constant? Generating examples for this task from the data is difficult. The particular minute a dosage is changed depends to a large extend on external conditions (e.g. an emergency involving a different patient). So interventions can be delayed and the optimal minute an intervention should be performed is unknown. To make sure that we generate examples only when a physician was closely monitoring the patient, we consider only those minutes where some drug was changed. This leads to 1319 training and 473 test examples.

Table 4. Confusion matrix for predicting time and direction of Dobutamine and Adrenaline interventions

	actual intervention		
	up	equal	down
Dobutamine			
predicted up	46	32	3
predicted equal	50	197	54
predicted down	5	30	56

	actual intervention		
	up	equal	down
Adrenaline			
predicted up	23	22	3
predicted equal	21	310	15
predicted down	4	34	41

For each drug we trained two binary SVMs. One is trained on the problem “increase dosage” vs. “do not increase dosage (i.e. lower or keep dosage equal)”, the other one is trained on the problem “lower dosage” vs. “do not lower dosage (i.e. increase or keep dosage equal)”. An intervention is predicted if exactly one such decision rule recommends a change. As an example, Figure 2 shows the decision rule that the SVM learned for increasing the dosage of Glyceroltrinitrate. Since the class distribution is very skewed towards the “do not ... dosage” class, we use a cost model. The cost-factors are chosen so that the potential total cost of the false positives equals the potential total cost of the false negatives. This means that the parameters of the SVM are chosen to conform to the ratio

$$\frac{C_+}{C_-} = \frac{\text{number of negative training examples}}{\text{number of positive training examples}} \quad (4)$$

Table 4 shows the test results for Dobutamine and Adrenaline. The confusion matrices give insight into the class distributions and the type of errors that occur. The diagonal contains the test cases, where the prediction of the SVM was the same as the actual intervention of the physician. This accounts for 63% of the test cases for Dobutamine and for 79% of the test cases for Adrenaline. The SVM suggests the opposite intervention in about 1.5% for both drugs.

Table 5. Confusion matrix for predicting time and direction of Dobutamine and Adrenaline interventions in comparison to human performance (results from an experienced intensivist in brackets).

Dobutamine	actual intervention		
	up	equal	down
predicted up	10 (9)	12 (8)	0 (1)
predicted equal	7 (9)	35 (31)	9 (9)
predicted down	2 (1)	7 (15)	13 (12)

Adrenaline	actual intervention		
	up	equal	down
predicted up	4 (2)	3 (1)	0 (0)
predicted equal	4 (6)	65 (66)	2 (2)
predicted down	1 (1)	8 (9)	8 (8)

Again, we would like to put these numbers into relation to the performance of an expert when given the same information. For a subsample of 95 examples from the test set, we asked a physician to perform the same task as the SVM. The results for Dobutamine and Adrenaline are given in Table 5. The results of the SVM on this subsample are followed by the performance of the human expert in brackets. Both are aligned remarkably well. Again, the learned functions of the SVM are comparable in terms of accuracy with a human expert. This also holds for the other drugs.

3.4 SVM Rules in Evidence Based Medicine

To use the SVM decision functions in a bigger learning environment the binary decisions of the SVM often do not offer enough information to decide for the appropriate action. For example a decision to increase a drug may have been triggered by random effects in the data or different decision rules may advise two or more contradicting actions. Hence, a measure of evidence of the SVM decisions would be very useful. The numerical value of the SVM function $f(\vec{o}) = \vec{w} \cdot \vec{o} + b$ (remember that the SVM decision function is given by $h(\vec{o}) = \text{sign}(f(\vec{o}))$) can be used as such a measure [23].

As an example, Figure 3 shows the actual dosage of adrenaline of a patient over a period of 110 minutes (upper line) and compares this to the output of the SVM that was trained to the task of classifying whether or not to increase the dose of adrenaline (lower line). It can be seen that the SVM did recommend to increase the dosage for some time before the intervention took place, but the evidence to intervene rapidly increases a few minutes before the actual intervention. Shortly after the intervention the recommendation of the SVM quickly changes to “do not increase the dosage”.

From the viewpoint of quality control in medicine, the question whether the intervention should have been taken some time earlier, as the output of the SVM indicates, deserves further investigation. This might be an example of a situation where a more sophisticated alarm system would have alerted the intensivist on duty much earlier.

3.5 More Learning Tasks

Let us now reason about the appropriate learning tasks for our goals. One may ask whether learning the appropriate direction of an intervention is justified at all, or whether the real task is to find the

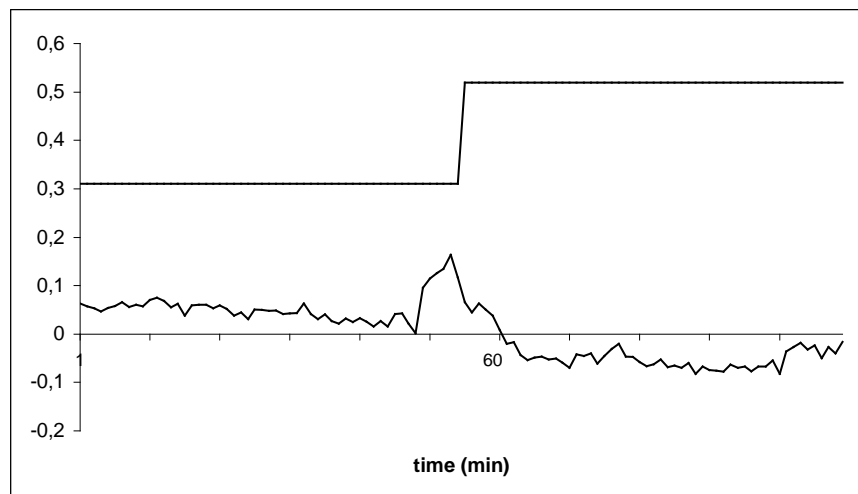


Figure 3. Actual dose of adrenaline (upper line) and evidence of SVM for “increase dosage” (lower line).

optimal dosage of a drug. In other words: should medical interventions be modeled as a classification or a regression problem? This is how we try to answer the question: For every drug, medical reasoning gives a value δ which a dosage change has to exceed in order to be considered to have a significant effect. We found that for all drugs at least 84% of the changes (96% at the average) lie within the range of $\pm \delta$. This justifies our approach. A higher dosage change can be realized by re-evaluating the decision to increase / decrease a drug a few minutes after the intervention.

Another interesting learning task would be to predict a trend in the vital signs of a patient. Discovering life-threatening situations as early as possible is a major key for optimal medical treatment. Moreover, this would also bring important advantages from the viewpoint of quality control and knowledge revision: In the validation of the effects of medical interventions, both of human experts and computer systems, one often finds cases where the expected effect of an intervention cannot be found in the data (e.g. medical knowledge says that the application of a drug will increase the blood pressure but the blood pressure stays stable). Confronted with this contradiction, a frequent explanation of experts is that the intervention anticipated an imminent change of the patients state into the opposite direction. As it is impossible to do a controlled experiment where the reaction of the patient with and without the intervention can be compared, the prediction of the patient's state based on examples of time periods without an intervention could offer a possibility to validate the success of an intervention.

Unfortunately, our experiments to predict vital signs of a patient in the nearer future (5 to 30 minutes) failed. For each vital sign, a regression version of the SVM [26] learned how much the parameter would increase or decrease. The learning results failed to predict these changes with more than default accuracy. As we tried many different representations of the patient's state (with and without history, learning an individual predictor per patient vs. learning on all patients, using Fourier transforms of the measurements of vital signs), we feel that this learning task is ill-posed. At the level of numerical measurements, i.e. disregarding the qualitative knowledge about physiological processes,

the prediction cannot become more precise. Hence, we combine data-driven numerical methods with a knowledge-based approach (see Section 4).

Another learning task could aim at characterizing a stable state by the observed measurements. Instead of judging the patients state in term of necessary interventions, a learning algorithm could find a description of regions in the high-dimensional attribute space that can be considered safe. When the patients state leaves the safe regions, an alarm can be generated. There exists an extension of the SVM algorithm to estimate the support of high-dimensional data [24] that seems to be promising for this learning tasks.

4 Medical Knowledge Base

Decision rules learned by the SVM reflect the average behavior of a physician, not the “gold standard”. As argued above, they have to be checked against medical knowledge about the effects of drugs. This section presents an approach to building a knowledge base that helps accomplish this task automatically and that makes decision support transparent.

Knowledge acquisition from experts is performed according to the current state of the art: first, knowledge is elicited from the expert, second, a knowledge base is modeled, third, the model is inspected, validated, and enhanced in collaboration with the expert. These steps form a cycle, i.e. the third step actually leads to obtain more expert knowledge, which is then modeled, etc.[19]. This expert knowledge augments and validates the data-driven knowledge acquisition using machine learning.

4.1 Knowledge Acquisition and Representation

The knowledge base of action-effect rules serves three purposes. First, it is used in order to model a protocol of care. Second, it is used to base learned decision functions on explicit and qualitative knowledge. Third, it is used for the validation of predictions. Let us describe the knowledge acquisition from experts before we show how this

Table 6. Medical Knowledge base for hemodynamic effects: + = increase of the respective variable or intervention; - = decrease; 0 = no change.

Intervention	Effect on hemodynamic variable					
	Heart Rate	Mean Arterial Pressure	Mean Pulmonary Artery Pressure	Central Venous Pressure	Cardiac Output	
Dobutamine	+	+	+	+	0	+
	-	-	-	-	0	-
Adrenaline	+	+	+	+	0	+
	-	-	-	-	0	-
Noradrenaline	+	-	+	+	0	-
	-	+	-	-	0	+
Nitroglycerin	+	+	-	-	-	+
	-	-	+	+	+	-
Fluid intake / output	+	-	+	+	+	+
	-	+	-	-	-	-

knowledge is integrated with the learned decision functions (section 4.3) and how it is used for validating predictions (section 5).

A medical expert defined the necessary knowledge. This knowledge is medical textbook knowledge for the cardiovascular system. It reflects direct pharmacological effects of a selected list of medical interventions on the basic hemodynamic variables. Any interaction of these interventions with other organ systems or of other organ systems with the cardiovascular system were ignored. An excerpt of intervention-effect relations is shown in Table 6. The dosage intervals indicated for each drug are not shown in the table, but modeled in the knowledge base. Also parameter dependencies have been modeled. It should be noted that the knowledge is qualitative with intervals of dosages, trends of changes, and implicit time intervals.

For the representation of qualitative medical knowledge we chose the MOBAL system [20]. MOBAL is a knowledge acquisition and maintenance system. Several tools facilitate the construction and inspection of a knowledge base. Its representation formalism is a restricted many-sorted first-order logic with explicit negation. A four-valued logic is used in order to allow for unknown and contradictory facts in addition to true and false facts. The inference engine derives

new facts on the basis of rules and given facts. Due to the expressive power of first-order logic, compact models can be built. What would be a rule in propositional logic, can be expressed by a mere fact in first-order logic. For instance, using a propositional logic, explicitly stating that up is the opposite of down requires the rule

heart_rate_trend=up --> not (heart_rate_trend=down)

and its dual form for all parameters. Using first-order logic, the fact

opposite(up, down)

is stated and can be used for any parameter. The pharmacological knowledge from Table 6 is expressed by facts of the form

effect(adrenaline, 0.01, 0.03, art, up)

stating that Adrenaline in a dosage between 0.01 and 0.03 mg/kg/min has the effect up on mean arterial pressure. Effects are modeled for substances. Additional facts indicate the particular drugs in which the substance is contained.

Patient records are also expressed by facts. The time is indicated by minutes, starting with the first measurement of a patient and ending with his or her discharge from intensive care.

intervention(pat4711, 10, 62, supra,0.02)

means that the patient 4711 from the tenth minute to minute 62 received Suprarenin (a drug containing Adrenaline) in a dosage of 0.02 mg/kg/min. Given the abstractions described in section 2, the values of hemodynamic parameters are stated in terms of level changes.

level(pat4711, 11, 62, hr, up)

states that the heart rate of patient 4711 had an upward level change at minute 11 and then remained almost stable until minute 62. In addition to this abstract description of a vital sign in a time interval, its deviation from the stable state is calculated. For each vital sign, the desired range of values is given, e.g. [60, 100] for the heart rate. For a patient's parameter values within a time interval, the standard deviation is calculated and added to (subtracted from) the upper (lower) value of the desired range. If the patient's actual value does not lie within this enlarged interval, a fact stating a deviation is entered. For instance, the following fact states that arterial mean pressure of patient 4999 is beyond the desired range:

deviation(pat4999, 0, 31, art, up)

We now want to use the pharmacological knowledge for deriving expected effects of an intervention on a particular patient. This is done by rules. The advantage of first-order logic is particularly important for modeling relations between intervals. For instance, stating that two time intervals are immediately succeeding, can be expressed by simply unifying the end point of one time interval with the start point of the other time interval. The following statement states, for instance, that two interventions were directly succeeding each other:

intervention(Patient, T1, T2, M, D1)

intervention(Patient, T2, T3, M, D2)

This statement can be instantiated by all patients, points in time, parameters and dosages as long as the same argument variable (e.g. Patient) is instantiated by the same value (e.g., pat4711). Different argument variables (e.g. D1, D2) can be instantiated by different values.

intervention(pat4711, 73, 83, supra, 0.05)

intervention(pat4711, 83, 177, supra, 0.02)

Intervals of dosages are handled in a similar manner. We can distinguish between major and minor changes of a dosage. A minor change is one within the same interval for which an effect has been stated by pharmacological facts. The rule and an actual instantiation is the following:

intervention(Patient, T1, T2, M, D1),

intervention(Patient, T2, T3, M, D2),

contains(M, S),

effect(S, FromD1, ToD1, Param, Trend),

FromD1 = < D1 < ToD1, FromD1 = < D2 < ToD1

-->

interv_effect(Patient, T2, T3, M, Param, Trend, minor)

intervention(pat4711, 441, 968, nitro, 1.9),

intervention(pat4711, 968, 1081, nitro, 2.38),

contains(nitro, glyceroltrinitrat),

effect(glyceroltrinitrat, 1, 10, hr, up),

1 = < 1.9 < 10, 1 = < 2.38 < 10

-->

interv_effect(pat4711, 968, 1081, nitro, hr, up, minor)

Changing into another such interval is a major change. The actual dosage of a drug given to a patient is compared with the dosage interval

of effect facts. The following rule expresses the enforcement of an effect because of a major change of dosage.

```
intervention(Patient, T1, T2, M,D1),
intervention(Patient, T2, T3, M,D2),
contains(M, S),
effect(S, FromD1, ToD1, Param, Trend),
effect(S, FromD2, ToD2, Param, Trend),
FromD1 =< D1 < ToD1, FromD2 =< D2 < ToD2,
ToD1 < FromD2
-->
interv_effect(Patient, T2,T3, M, Param, Trend, major)
```

Note, that if the substance S of drug M has a decreasing effect on a parameter of the patient, the rule predicts a further decrease of that vital sign. The variable Trend is then instantiated by down. Another rule states that decreasing a substance with an increasing effect on a parameter will decrease the parameter's value. We use such rules in order to predict effects of interventions. The prediction of intervention effects is used to check interventions that are proposed by the learned decision rules. Not counting the patient records, the knowledge base consists of 39 rules and 88 facts.

4.2 Validating Action-Effect Rules

In order to validate the knowledge base we applied it to the data of 148 patients. The data contain 8,200 interventions. The validation is easy, since rules can directly be applied to patient data. MOBAL's inference engine derived 27,400 effects of the interventions using forward chaining. For 22,599 effects the actual effects in terms of level changes could be computed by the time series analysis (see section 2). When matching the derived effects with the actual ones, the system detected:

- 13,364 effects (i.e. 59.14%) took place in the restricted sense, that the patient's state remained stable. E.g., a drug with an increasing effect on a patient's vital sign does not lead to a significant level change of this parameter. This is not in conflict with medical knowledge, but shows best therapeutical practice. Smooth medication keeps the patient's state stable and does not lead to oscillating reactions of the patient.

- 5,165 effects (i.e. 22.85%) took place in the sense, that increasing or decreasing effects of drugs on vital signs match corresponding level changes.
- 4,070 contradictions (i.e. 18.01%) were detected. The observed level change of a vital sign went into the opposite direction of the knowledge-based prediction.

The ratio of 83.56 percent correct predictions of effects is quite positive. Some decisive features are not present in the data. Particularly the lack of data about cardiac arrhythmias and cardiac output could possibly explain many deviations of observed from predicted effects.

4.3 Integrating Learned Decision Functions With the Knowledge Base

Since the goal of our work is an integrated system for intensive care monitoring, the numerical approach using the SVM has to be incorporated into the logic of MOBAL. While training SVM classifiers can take place offline in a separate program, MOBAL needs to be able to evaluate SVM decision rules and access the results online. We achieve this by introducing the special predicate `svm_calc/6` with the following semantic. The first two arguments indicate the patient and the drug. The third argument is either “up” or “down” depending on whether the `svm_calc` fact belongs to the SVM predicting dose increase or decrease (compare section 3.3). The fourth argument is the time and the fifth is the current dosage of the drug. The last argument finally contains the value of that particular SVM rule for the measurements at that time. Calculating can be done very efficiently, since it mainly consists of computing a dot product between the SVM weight vector and the measurement vector. From each pair of decision rules (i. e. up and down) an intervention for the respective drug is recommended, if exactly one decision rule has a value larger than a confidence threshold of 0.8.

The decision rule for an increase of Glyceroltrinitrat (nitro) together with the actual parameter values of patient 4999 at time 32 is shown in Figure 2. The dot product plus -4.368 (the value of b) is 1.85598 . The fact entered into the fact base for patient 4999 is `svm_calc(pat4999,`

nitro, up, 32, 0.0, 1.85598). An intervention to increase nitro is derived. The dose is calculated on the basis of the former dose. The SVM actually only decides whether to increase, to decrease, or not to change the dose. For each drug, a level of granularity is defined. For instance, the granularity of Glyceroltrinitrat is 1, whereas that of Suprarenin (containing adrenaline) is 0.01. The dose is changed by just one step. In our example, the proposed intervention is:

pred_intervention(pat4999, 32, nitro, 1.0).

5 Using the Knowledge Base of Effects to Validate Interventions

Medical knowledge is used for validation in two different ways. On the one hand, learned decision rules are validated on patient data by comparing the effects of their recommended interventions with the effects of actual physicians' interventions. This validation means to incorporate an evaluation step already into the knowledge acquisition phase. On the other hand, we believe that even an evaluated decision support system should check its decisions by considering their effects.

5.1 Validating Learned Decision Rules

There are usually several different combinations of drugs that achieve the same goal of keeping the patient in a stable state. And indeed, different physicians, depending on their experience in the ICU, do use different mixtures and follow different strategies to reach this goal. For comparing treatment strategies, the real criterion is whether the recommendations have the same effect as the actual interventions. Therefore, we apply the action-effect rules from the knowledge base to both the proposed intervention of the SVM classifiers and to the intervention actually performed by the physician. If the derived effects are equal, then the proposed decision of the SVM classifiers can be considered as "equivalent" to the intervention executed by the physician. The results of this comparison for 473 interventions are shown in Table 7. The right-most column indicates the accuracy, i.e. in how many cases the classification of SVM and physician were identical (same behavior of SVM and physician). The other columns state how

Table 7. Equivalence of decisions regarding effects.

Interventions	Mean arterial pressure	Heart rate	Same effect all parameters	Same behavior
Dobutamine	403	395	383	299
Adrenaline	407	406	393	374
Glyceroltrinitrate	437	388	380	342
Noradrenaline	436	428	424	420
Nifedipine	457	457	455	438

often the SVMs' intervention leads to the same effects as the intervention of the physician. The first two columns show, how many of interventions had the same effect on arterial blood pressure or heart rate, respectively. The third column gives a more concise evaluation. Here it is stated, how many interventions recommended by the SVM had the same effects on all vital signs as the actual intervention. For instance, the SVM correctly classifies 299 test cases for Dobutamine (63%). If we compare the resulting effects of the predicted interventions concerning Dobutamine with the effects of the actual physician's interventions, we find that in 383 cases (81%) the deduced effects will be equal. Thus, in 84 cases the recommendation of the SVM does not match the physician's behavior, but the derived effects are the same, since the physician has chosen an "equivalent" drug or combination of drugs. An inspection of these cases helps to clarify issues of best practice and thus supports knowledge acquisition.

5.2 Validating Proposed Interventions

As depicted in the overall architecture (cf. Figure 1), we have chosen a design which allows us to use the action-effect rules in the knowledge base for validating predicted interventions. The underlying argument is that accuracy measures only reflect how well the SVMs' learning results fit actual behavior of the physician. However, we aim at best practice. Hence, we validate a proposed intervention with respect to its effects on the patient. If the effects push vital signs in the direction of the desired value range, the recommendation is considered sound, otherwise it is rejected. An example may clarify this. Patient 4999 is older than 75 years and stays at the ICU after a surgical operation. He suffers from high arterial mean pressure (around 124), where the heart rate is normal (around 80). Using its decision rules, the SVM

recommends to increase Glyceroltrinitrat (see Figure 2). This proposed intervention is checked by the medical knowledge about effects. The derived effects are an increase of the heart rate and a decrease of arterial mean pressure as well as left ventricular stroke work index (lvswi) and systemic vascular resistance (svr): *interv_effect(pat4999,32, T, art, down)*. The observed deviation is *deviation(pat4999, 0,31, art, up)*. Since down is the opposite of up, the proposed intervention is considered sound. In this way, the prescriptive medical knowledge (action-effect rules) is used to control the knowledge that is learned from actual therapies (state-action rules).

6 Comparison With Related Work

Using data from the most comprehensive singular clinical data repository at the LDS Hospital, Salt Lake City, Utah, USA, the group of Morris [22] developed a rule-based decision support system (DSS) for respiratory care in acute respiratory distress syndrome. Time is handled by introducing time points into the rules where a certain parameter value needs to be obtained. The development of this highly specialized system required more than 25 person years. It is a propositional rule base without a mechanism for consistency checking or matching rules and data. All validation efforts started only after the knowledge base had been completed.

Temporal reasoning is taken seriously in other developments [3, 5, 17, 25]. The Stanford approach uses an explicit time ontology for low-frequency data [25]. This approach is not feasible for our application. The VIE-VENT system is comparable with our approach in that it combines numerical data and a knowledge base [17]. Qualitative abstractions are derived for deviations of measurements from the target range. Time intervals refer to the validity of a measurement. The detection of outliers (data validation) is handled by a trend-based component. The validated measurements are used by the therapy planning component which aims at pushing vital signs into the value ranges of a stable state. Similar to our approach, therapy planning is divided into state-action rules (therapeutic actions based on status interpretation) and verifying the effectiveness of interventions. However, the system was developed without using actual patient data.

Hence, the observation that parameter values oscillate considerably was made as late as the first clinical experience. In contrast, this observation has motivated our use of the phase space procedure for abstracting from numerical time series. Temporal correlations can also be included in trend templates, which are used by Haimowitz and Kohane [5]. Trend templates consist of sets of low order polynomial regression models describing qualitative characteristics. Pattern abstraction is done based on the fit of these templates to the observed data. The major drawbacks of this method are the demand for predefined expected behavior and absolute value thresholds. However, time series in intensive care often show irregular behavior like patchy outliers, or outliers and level changes occurring in short time lags. Such behavior is difficult to specify in advance. Moreover, thresholds should be depending dynamically on the patient's status in the past. This has already been included in our approach, which does not need prespecified patterns either. Altogether, statistical time series analysis seems to be the most sophisticated method to model and investigate dynamical data since other approaches capture only parts of the time dependent structure of the data.

Our goals of easing the development of guidelines and validating the knowledge early on is shared by the two-step approach by Mani and coworkers [16]. They use machine learning in order to first characterize scores of dementia with respect to six categories (e.g., memory, orientation). These learning results are then used to learn the global clinical dementia rating. After a two years effort an efficient and effective system was accomplished. While the goals are the same, the application characteristics and, hence, the methods are completely different. The clinical rating is a classification task and the patient data is of qualitative nature, whereas our task is on-line monitoring and the patient data are time series of numerical measurements.

7 Conclusions

We presented an approach towards integrating learning and knowledge-based methods for the development of decision support algorithms in critical care. The SVM was chosen for learning state-action rules due to its ability to handle multiple features. For modeling medical knowledge

in terms of action-effect rules we chose a first-order logic representation using MOBAL. This allowed a compact representation of medical knowledge with a small number of rules, fulfilling the real-world demand for a knowledge base to be understandable by humans and accessible for expert validation.

The validation issue has been treated with special care. Each process has been validated in the standard way, i.e. tested on data not used for training. In addition, the results of state-action rules were compared with the results of a human expert who classified the same data. Moreover, recommended interventions of state-action rules are validated by formalized medical knowledge. On the one hand, the effect of a recommended intervention is compared with the effect of an actual intervention. Of course, this comparison can only be made for past cases. In case of conflict, the expert inspects the particular cases. This may lead to the generation of explicit additional knowledge. On the other hand, the formalized effects of interventions are applied to current cases and evaluated with respect to the target ranges of vital signs.

Our new approach combines modeling of expert knowledge with data-driven methods. This eases the task of building operational protocols. Moreover, the data-driven method allows for an ongoing enhancement of the knowledge base on the basis of current practice. The knowledge base is validated against existing patient data. This approach is meant to be significantly more effective than the tedious, time-consuming, and costly process of traditional development of on-line operational decision support systems. The effect of this is an improvement in both the extensibility of an existing knowledge base and the control of the quality of medical treatment.

8 Acknowledgements

This work has been funded in part by the Deutsche Forschungsgemeinschaft (SFB475, "Reduction of Complexity for Multivariate Data Structures").

9 References

- [1] Bauer, M., Gather, U., and Imhoff, M. (1999) "The identification of multiple outliers in online monitoring data," *Technical Report 29, SFB 475*, University of Dortmund.
- [2] Box, G. E. P., Jenkins, G. M., and Reinsel G. C. (1994), "Time Series Analysis. Forecasting and Control,". Third Edition, Prentice Hall, Englewood Cliffs.
- [3] Dojat, M. and Sayettat, C. (1995) "A Realistic Model for Temporal Reasoning in Real-Time Patient Monitoring", *Applied Artificial Intelligence*, Vol. 10 (2), pp.121-143.
- [4] Guyatt, G., Drummund, M., Feeny, D., Tugwell, P., Stoddart, G., Haynes, R., Bennett, K., and LaBelle, R.(1986) "Guidelines for the clinical and economic evaluation of health care technologies", *Soc Sci Med*, Vol. 22, pp.393-408.
- [5] Haimowitz, I. J., and Kohane, I. S. (1996) "Managing temporal worlds for medical trend diagnosis," *Artificial Intelligence in Medicine*, Vol. 8, pp. 299-321.
- [6] Imhoff, M. (1995) "A clinical information system on the intensive care unit: dream or night mare?," *Medicina Intensiva 1995, XXX. Congreso SEMIUC*. Murcia, pp. 17-22.
- [7] Imhoff, M. (1996), "3 years clinical use of the Siemens Emtex System 2000: Efforts and Benefits", *Clinical Intensive Care 7 (Suppl.)*, pp. 43-44.
- [8] Imhoff, M. (1998) "Clinical Data Acquisition: What and how?," *Journal für Anästhesie und Intensivmedizin*, Vol 5, pp. 85-86.
- [9] Imhoff, M., Bauer, M., Gather,U., and Löhlein, D. (1998) "Statistical pattern detection in univariate time series of intensive care on-line monitoring data," *Intensive Care Med*, Vol 24 pp. 1305-1314.

- [10] Imhoff, M., Bauer, M. and Gather, U. (1999) "Time-effect relations of medical interventions in a clinical information system" *KI-99: Advances in Artificial Intelligence*, LNAI 1701, Springer-Verlag, pp. 307-310.
- [11] Imhoff, M., Lehner, J.H. and Löhlein, D. (1994) "2 years clinical experience with a clinical information system on a surgical ICU," *7th European Congress on Intensive Care Medicine*, pp. 163-166.
- [12] Jennings, D., Amabile, T. and Ross, L. (1982) "Informal covariation assessments: Data-based versus theory-based judgements" *Judgement under uncertainty: Heuristics and biases*, Cambridge University Press, Cambridge, pp. 211-230.
- [13] Joachims, T. (1999) "Making Large-Scale SVM Learning Practical," *Advances in Kernel Methods - Support Vector Learning*, MIT Press, Cambridge, pp. 169-184.
- [14] Keravnou, E. T. (1996) "Temporal diagnostic reasoning based on time-objects," *Artificial Intelligence in Medicine*, Vol 8, pp. 235-265.
- [15] Kukar, M., Kononenko, I., Groselj, C., Kralj, K., and Fettich, J. (1999) "Analysing and Improving the Diagnosis of Ischaemic Heart Disease with Machine Learning," *Artificial Intelligence in Medicine*, Vol 16, pp. 25-50.
- [16] Mani, S., Shankle, W.R., Dick, M. B., and Pazzani, M. J. (1999) "Two-Stage Machine Learning for Guideline Development," *Artificial Intelligence in Medicine*, Vol 16, pp. 51-71.
- [17] Miksch, S., Horn, W., Popow, C., and Paky, F. (1996) "Utilizing Temporal Abstraction for Data Validation and Therapy Planning for Artificially Ventilated Newborn Infants," *Artificial Intelligence in Medicine*, Vol 8, pp. 543-576.

- [18] Miller, G. (1956) "The magical number seven, plus or minus two: Some limits to our capacity for processing information," *Psychol Rev*, Vol 63, pp. 81-97.
- [19] Morik, K.(1994) "Balanced Cooperative Modeling," *Machine Learning - A Multistrategy Approach*, Morgan Kaufmann, pp. 295-318.
- [20] Morik, K., Wrobel, S., Kietz, J.-U., and Emde, W. (1993) *Knowledge Acquisition and Machine Learning - Theory, Methods, and Applications*, Academic Press, London.
- [21] Morris, A. and Gardner, R. (1992) "Computer applications," *Principles of Critical Care*, McGraw-Hill, New York, pp. 500-514.
- [22] Morris, A. (1998) "Algorithm-Based Decision-Making," *Principles and Practice of Intensive Care Monitoring*, McGraw-Hill, New York, pp. 1355-1381.
- [23] Platt, C. (1999) "Probabilistic Outputs for Support Vector Machines and Comparisons to Regularized Likelihood Methods," *Advances in Large Margin Classifiers*, MIT Press.
- [24] Schölkopf, B. and Williamson, R. and Smola, A. and Shawe-Taylor, J. (1999) "SV-Estimation of a Distribution's Support", NIPS 99.
- [25] Shahar, Y. and Musen, M. (1996) "Knowledge-Based Temporal Abstraction in Clinical Domains," *Artificial Intelligence in Medicine*, Vol. 8, pp. 267-298
- [26] Vapnik, V. (1998) *Statistical Learning Theory*, Wiley, New York.