

Implementing an Integrated Time-Series Data Mining Environment - A Case Study of Medical KDD on Chronic Hepatitis -

H. Abe¹ and T. Yamaguchi²

¹Department of Medical Informatics, Shimane University, Shimane, Japan

²Department of Administration Engineering, Keio University, Japan

Abstract – In this paper, we present the implementation of an integrated time-series data mining environment. Time-series data mining is one of key issues to get useful knowledge from data bases. To execute time-series data mining smoothly, we have designed an environment which integrates time-series pattern extraction methods, rule induction methods and rule evaluation methods with active human-system interaction. After implementing this environment, we have done a case study to mine time-series rules from blood/urine biochemical test data base on chronic hepatitis patients. Then a physician has evaluated and refined his hypothesis on this environment. We discuss the availability of how much support to mine interesting knowledge for an expert.

Keywords – Data Mining, Time-Series Data, Time-Series Pattern Extraction, Rule Induction

I. INTRODUCTION

In recent years, EBM (Evidence Based Medicine) has been widely recognized as a new medical topic to care each patient with the conscientious, explicit and judicious use of current best evidence, integrating individual clinical expertise with the best available external clinical evidence from systematic research and the patient's unique values and circumstances. These evidences have been often found out in clinical test databases, which are stored on HIS (Hospital Information Systems). Looking at such findings, time-series rules are one kind of important medical evidences related to clinical courses of patients. However, it is difficult to find out such evidences systematically. Medical researchers need some systematic method to find out these evidences faster.

Besides, KDD (Knowledge Discovery in Databases)[3] has been known as a process to extract useful knowledge from databases. In the research field of KDD, Time-Series Data Mining is one of important issues to mine useful knowledge such as patterns, rules, and structured descriptions for a domain expert. Although time-series data mining can find out useful knowledge in given data, it is difficult to find out such knowledge without cooperation among data miner, system developers, and domain experts.

To above problems, we have developed a time-series data mining environment, which can apply medical data mining to find out medical evidences systematically, considering cooperation among data miners, system developers, and medical experts. Through a case study with a chronic hepatitis database, we have identified the

procedures, which have been needed to execute time-series data mining cooperatively with active human-system interaction. With this analysis, we have developed a time-series data mining environment, integrating time-series pattern extraction, rule induction, and rule evaluation through visualization/evaluation/operation interfaces.

In this paper, we present an implementation of the integrated time-series data mining environment with a case study of time-series rule mining on chronic hepatitis dataset.

II. RELATED WORK

Many efforts have been done to analyze time-series data at the field of pattern recognitions. Statistical methods such as ARIMA and autoregressive model have been developed to analyze time-series data, which have equalized sampling rate, linearity, and periodicity. As signal processing methods, Fourier transform, Wavelet, and fractal analysis method have been also developed to analyze such well formed time-series data. These methods based on mathematic models restrict input data, which are well sampled. However, time-series data include ill-formed data such as clinical test data of chronic disease patients, purchase data of identified customers, and financial data based on social events. To analyze these ill-formed time-series data, we take another time-series data analysis method such as DTW (Dynamic Time Wrapping)[1], time-series clustering with multiscale matching[5], and finding Motif based on PAA (Piecewise Approximation Aggregation)[6].

As the one of the methods to find out useful knowledge depended on time-series, time-series/temporal rule induction method such as Das's framework[2] have been developed. We can extract time-series rules in which representative patterns are expressed as closes of their antecedent and consequent with this method.

To succeed in a KDD process, human-system interaction especially need at data pre-processing and post-processing of mined result. Current time-series data mining frameworks focus on not human-system interaction but automatic processing. We introduce human-system interaction for data pre-processing and post-processing of mined result with visualization and evaluation interfaces for patterns and if-then rules based on time-series patterns.

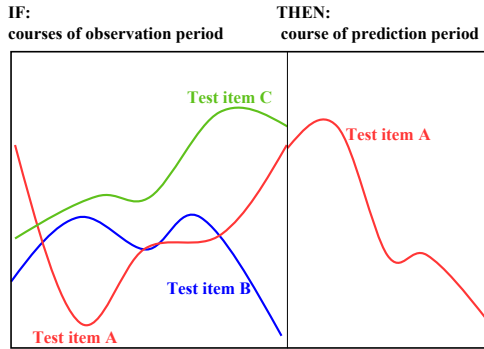


Fig. 1. Typical output if-then rule of our integrated time-series data mining environment.

With this environment, we aim the following efforts for each agent:

- 1) Developing and improving time-series data mining procedures for system developers
- 2) Collaborative data processing and rule induction for data miners
- 3) Active evaluation and interaction for domain experts

To implement the environment, we have analyzed time-series data mining frameworks. Then we have identified procedures for pattern extraction as data pre-processing, rule induction as mining, and evaluation of rules with visualized rule as post-processing of mined result. The system provides these procedures as commands for users. At the same time, we have designed graphical interfaces, which include data processing, validation for patterns on elemental sequences, and rule visualization as graphs. Fig. 2 shows us a typical system flow of this time-series data mining environment.

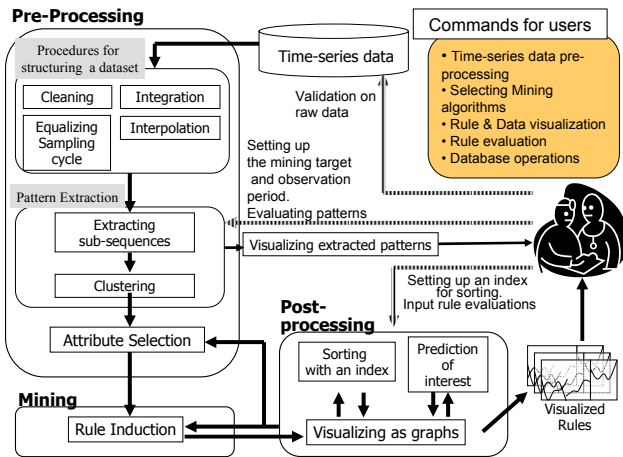


Fig. 2. An overview of the integrated time-series data mining environment with system flow.

A. Procedures to mine time-series rules

We have identified procedures for time-series data mining as follows:

- Data pre-processing
- pre-processing for data construction

- time-series pattern extraction
- attribute selection

Mining

- rule induction

Post-processing of mined results

- visualizing mined rule
- rule selection
- supporting rule evaluation

Other database procedures

- selection with conditions
- join

As data pre-processing procedures, pre-processing for data construction procedures include data cleaning, equalizing sampling rate, interpolation, and filtering irrelevant data. Since these procedures are almost manual procedures, they strongly depend on given time-series data and a purpose of the mining process. Time-series pattern extraction procedures include determining the period of sub-sequences and finding representative sequences with a clustering algorithm. We have taken our original pattern extraction method[11] for the case study described in Section IV. Attribute selection procedures are done by selecting relevant attributes manually or using attribute selection algorithms[7].

At mining phase, we should choose a proper rule induction algorithm with some criterion. There are so many rule induction algorithms such as Version Space[9], AQ15[8], C4.5 rules[12], and any other algorithm. To support this choice, we have developed a tool to construct a proper mining application based on constructive meta-learning called CAMLET. However, we have taken PART[4] implemented in Weka[14] in the case study to evaluate improvement of our pattern extraction algorithm.

To validate mined rules correctly, users need readability and ease for understand about mined results. We have taken 39 objective rule evaluation indexes to select mined rules[10], visualizing and sorting them depended on users' interest. Although these two procedures are passive support from a viewpoint of the system, we have also identified active system reaction with prediction of user evaluation based on objective rule evaluation indexes and human evaluations.

Other database procedures are used to make target data for a data mining process.

Since the environment has been designed based on open architecture, these procedures have been able to develop separately. To connect each procedure, we have only defined input/output data format.

B. Designing Interfaces to Collaborate Medical Experts and Data Miners

We have designed user interfaces as shown in Fig. 3.

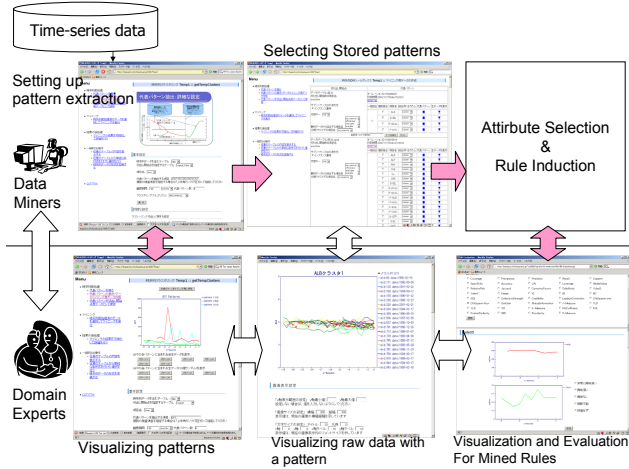


Fig. 3. Interfaces of the integrated time-series data mining environment.

For data miners, setting pattern extraction and selecting stored patterns to construct a dataset have been assigned. Pattern extraction setting interfaces are separated two phases. In first phase, a data miner set up basic settings such as an input time-series data set, a period for patterns, a target dataset for the purpose of the data mining process, and a clustering algorithm for pattern extraction. Then other detailed parameters such as the period for equalizing sampling rate and maximum interpolation period are set up at detailed setting phase. A data miner also set up detailed parameters of selected clustering algorithm at this phase. Besides, visualization of patterns, visualization of given raw data included selected pattern, and visualization of mined rules have been assigned for domain experts. Each visualization interface can be switched each other depended on their interest.

IV. A CASE STUDY WITH CHRONIC HEPATITIS DATABASE

As a case study of medical KDD, we have taken a clinical test dataset from Chiba University Hospital[13], which includes clinical blood and urine test data on chronic hepatitis B and C. Their test intervals are not only daily, weekly, or monthly but also randomized intervals depended on patients' statements. This dataset has approximately 1.6 million records. Each record consists of MID, test date, test item name, and test result. This dataset includes 771 patients, who have up to 20 years as his/her treatment period. We have taken IFN (interferon) treatment results of 195 patients as the target of rules, which represent as if-then classification rules. Each instance of this target data consists of MID, start date of IFN treatment, end date of it, his/her treatment result decided with GPT (ALT) values, and his/her treatment result based on virus markers. We have set up observation periods before finishing IFN treatment.

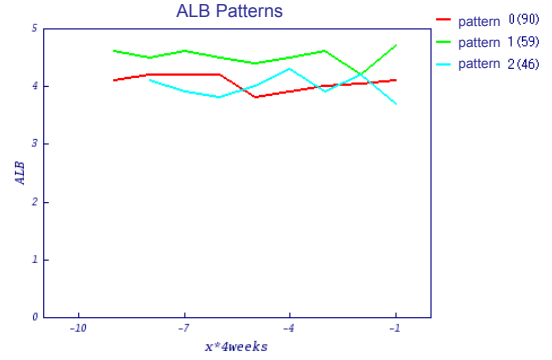


Fig. 4. Representative patterns on ALB while IFN treatment. '-1' week means finishing date of IFN treatment for each patient.

A. Phase1: Rousing new hypothesis in an expert

We presented extracted patterns about some observation periods to a physician. He noticed distinguishing patterns on ALB as shown 'pattern 1' and 'pattern 2' in Fig. 4.

He validated this pattern with each sequence of patients. He said that 'pattern 0' indicates typical course while in IFN treatment period. However, the other patterns show remarkable courses based on his knowledge. Since he thought 'pattern 2' shows what patients included in this pattern had adverse reactions at the end of IFN treatment period, he interested in their treatment results. On the other hand, 'pattern 1' indicates less reaction to IFN treatment at the end of the period. He roused a hypothesis that patients, who indicate good reaction after IFN treatment, have moderate adverse reaction during IFN treatment.

B. Phase2: Ensuring expert's hypothesis

To endure the hypothesis, we have extracted patterns of 0.8 years (40 weeks) as observation period for 40 test items. After joining these patterns as attributes of the dataset, we have induced if-then rules with the dataset. Then the patterns and rules are evaluated by the physician, visualizing patterns, rules and patient data on demand.

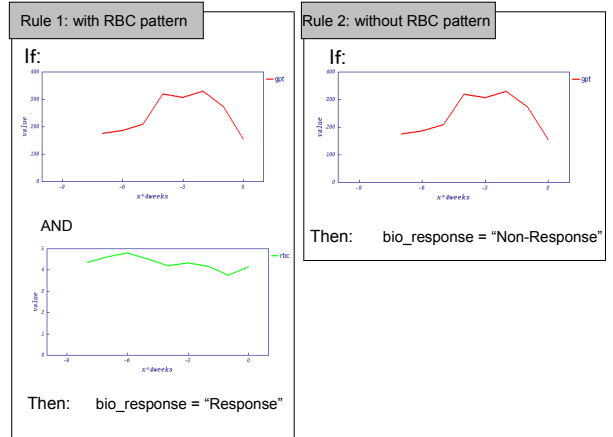


Fig. 5. Two rules, which have opposite class values, because of just one pattern at the antecedent of the left rule.

He interested in the rule (Fig. 5). The pattern of RBC shows trend of anemia, which is one of typical adverse reactions of IFN treatment. These patients included in this pattern would be anemia, he said.

He noticed that some combination of patterns show good result after IFN treatment but other combinations of patterns show no good result without or changing just one pattern.

At the same time, we have also developed pattern extraction method based on irregular sampling and quantization. Table I shows us numbers of expert's evaluations. The number of 'very interesting' increases after improvement. This result has been caused by one of the reasons why the improvement of pattern extraction algorithm extracts patterns, which can be suitable to his knowledge about patients' course.

TABLE I
NUMBERS OF EVALUATED RULES FROM BEFORE AND AFTER
IMPROVEMENT OF THE PATTERN EXTRACTION ALGORITHM

Evaluation Labels	Before improvement	After improvement
Very Interesting	4	15
Interesting	7	5
Fair	15	11
Difficult to understand	5	1
TOTAL	31	32

V. CONCLUSION

We have implemented a time-series data mining environment, which integrates time-series pattern extraction, rule induction, and rule evaluation with active interaction on graphical interfaces. As the result of the case study on a chronic hepatitis dataset, we have succeeded in rousing and ensuring hypothesis about IFN treatment in the mind of a medical expert.

We are connecting subsystems for selecting a proper attribute algorithm and a proper rule induction algorithm to this environment. Then we will construct time-series data mining application based on active user reaction, acquiring with rule evaluation interface.

ACKNOWLEDGMENT

We would like to thank Dr. Hideto Yokoi, who is the vice manager of Department of Medical Informatics, Kagawa University Hospital. He gave us useful comments and valuable evaluations about chronic hepatitis data mining during whole process of development of the integrated time-series data mining environment.

We also thank Dr. Miho Ohsaki, who is in Faculty of Engineering, Doshisha University. She gave valuable ideas

and precious discussions for the post-processing interface design of our time-series data mining system.

REFERENCES

- [1] D. J. Berndt and J. Clifford, "Using dynamic time wrapping to find patterns in time series", in *Proc. of AAAI Workshop on Knowledge Discovery in Databases*, 1994, pp.359-370.
- [2] G. Das, L. King-Ip, M. Heikki, G. Renganathan, and P. Smyth, "Rule Discovery from Time Series", in *Proc. of International Conference on Knowledge Discovery and Data Mining*, 1998, pp.16-22.
- [3] U. M. Fayyad, G. Piatctsky-Shapiro, and P. Smyth, "From Data Mining to Knowledge Discovery: An Overview", *Advances in Knowledge Discovery and Data Mining*, AAAI Press / The MIT Press, CA, 1996, pp.1-34.
- [4] E. Frank, I. H. Witten, "Generating accurate rule sets without global optimization", in *Proc. of the Fifteenth International Conference on Machine Learning*, 1998, pp.144-151.
- [5] S. Hirano and S. Tsumoto, "Mining Similar Temporal Patterns in Long Time-Series Data and Its Application to Medicine", in *Proc. of the 2002 IEEE International Conference on Data Mining*, 2002, pp.219-226.
- [6] J. Lin, E. Keogh, S. Lonardi, and P. Patel, "Finding Motifs in Time Series", in *Proc. of Workshop on Temporal Data Mining*, 2002, pp.53-68.
- [7] H. Liu and H. Motoda, "Feature selection for knowledge discovery and data mining", Kluwer Academic Publishers, 1998.
- [8] R. Michalski, I. Mozetic, J. Hong, and N. Lavrac., "The AQ15 Inductive Learning System: An Overview and Experiments", *Reports of Machine Learning and Inference Laboratory*, MLI-86-6, George Mason University, 1986.
- [9] T. M. Mitchell, "Generalization as Search", *Artificial Intelligence*, 18(2), 1982, pp.203-226.
- [10] M. Ohsaki, S. Kitaguchi, K. Okamoto, H. Yokoi, and T. Yamaguchi, "Evaluation of Rule Interestingness Measures with a Clinical Dataset on Hepatitis", in *Proc. of ECML/PKDD 2004*, LNAI3202, 2004, pp.362-373.
- [11] M. Ohsaki, H. Abe, S. Kitaguchi, S. Kume, H. Yokoi, and T. Yamaguchi, "Development and Evaluation of an Integrated Time-Series KDD Environment - A Case Study of Medical KDD on Hepatitis-", *Joint Workshop of Vietnamese Society of Artificial Intelligence, SIGKBS-JSAI, ICS-IPSI and IEICE-SIGAI on Active Mining*, 2004, No.23.
- [12] J. R. Quinlan, "Programs for Machine Learning", Morgan Kaufmann, 1992.
- [13] S. Tsumoto, Hepatitis Dataset for Discovery Challenge, <http://lisp.vse.cz/challenge/ecmlpkdd2002/index.html>, 2002.
- [14] I. H. Witten and E. Frank, "Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations", Morgan Kaufmann, San Francisco, 2000.