

Information Summarization for Network Performance Management

Mikko Kylväjä^N, Pekka Kumpulainen^T, Kimmo Hätönen^{NR}

(N) Nokia Networks, Espoo, Finland

(T) Tampere University of Technology, Tampere, Finland

(NR) Nokia Research Center, Helsinki, Finland

Abstract – Mobile communication networks have grown and the supporting technology is increasingly complicated. Operating a cellular mobile network system is challenging and requires expert effort. The requirements to improve efficiency and to reduce the operating expenses demand the cellular network operators to apply automated solutions in performance management of the network systems. Experts at operators' network management centre must analyze large amounts of data. They spend most of their time doing routine tasks, such as data acquisition, filtering and repetitive decisions. This paper presents an information summarization method to support the experts in decision-making and to reduce the effort needed to analyze large amounts of data. Furthermore, we expect the method to improve the decisions as human errors are reduced. The method was tested by applying it to a data set from a commercial GSM mobile network.

Keywords: telecommunication, mobile network, performance management

1. INTRODUCTION

The analysis method presented here is aiming to support experts in their routine performance management of GSM networks. Operating a mobile network involves a series of complicated tasks [1]. The performance management expert must decide on operative actions on element or network parameters so that the performance will be satisfactory for the end users. At present, the expert must laboriously collect large amounts of data. Once having the data the challenge is that it is excessive for human understanding. An expert needs to find manually the relevant part from the data for present decision-making, and that task is repeated several times during a day. All this is quite inefficient use of the experts' time without tools that automate the routine tasks and focus their attention to the most relevant part of the data.

A mobile network produces 3-dimensional data matrices shown in Fig. 1. Several indicators are stored for a number of network elements at a predefined sampling interval.

For performance management purposes the data is usually summarized in time so that hourly or daily averages are used. Each indicator is summarized from several original measurement variables with a predefined formula to Key Performance Indicator (KPI) [2].

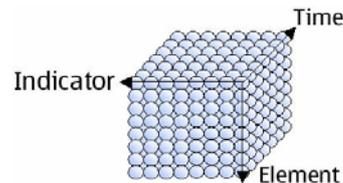


Fig. 1. Summarization of information can be done in three dimensions.

Our method provides summarization in element dimension by presenting behaviour of problem classes rather than individual elements. The indicators are further summarized by pointing out the ones that most significantly distinguish the element's class from other problem classes.

Section 2 is a survey to research and related work done in this area. Section 3 describes the decision support system that is used to reduce the data and to process it to help the expert in decision-making.

In section 4 we present a case study and results of real network scenario. The processed data was shown to radio network expert and his analysis of the results is presented.

2. RELATED WORK

Data reduction strategies and concepts [3] have also been applied in network management:

- *Data aggregation*; it is common in network management that performance data is aggregated for parent objects in network hierarchy. Another common aggregation for performance data is summarization to longer time-periods.
- *Dimension reduction*; because the amount of available network performance measurements is large (several hundreds) it is important to select only the ones that are relevant for the goal of the expert. In network management systems there are typically several ready-made reports that contain the indicators for only one network functionality (for example, GRPS or capacity).
- *Data compression*; data compression is used especially in radio interface, where the amount of bits used to express a message is minimized.
- *Numerosity reduction*; Clustering similarly behaving elements to groups is a typical example to reduce the numerosity.

- *Discretization and hierarchy generation*; an example of hierarchical data reduction is definition of high-level performance indicators that is a combination of lower level indicators.

There are several industrial application areas where information summarization is achieved by clustering the elements or samples according to behaviour. Clustering enhances efficiency in processes that are complex and provide data sets with numerous high dimensional samples. Paper mill is an example, where process quality has been analysed with Self-organising maps and other clustering methods [4, 5].

In telecommunications, clustering has been applied to measurement data [6]. Another important application area is alarm clustering, which is efficient in intrusion and fraud detection. [7] Also correlating alarms in a GSM network management system has been successfully applied in commercial networks [8, 9]. In all cases the target is to improve the efficiency by reducing the amount of alarms or measurements histories shown to operator by summarizing them according to common root cause.

Clustering algorithm efficiency has been studied in order to scale up the sizes of the analyzed databases. [10] The efficiency can be obtained by optimizing the algorithm to handle higher amount of dimensions [11].

Clustering methods are suitable for customer behaviour analysis, which reveals useful information for marketing based on the purchasing patterns. [3] Customer behaviour analysis is also used in telecommunications to analyse the satisfaction of the users and to improve the dimensioning of the network and to optimize the pricing policy [18].

Other recent research areas, in which the clustering has been applied, are Machine Learning, Pattern Recognition and genetics, just to name a few [3, 12].

The results of the clustering are heavily dependent on the pre-processing and scaling of the input data. By utilising different scaling methods, different information is highlighted in the data set [13, 14].

3. DECISION SUPPORT SYSTEM

We present a method to automate the classification of network elements according to their state and to summarize the relevant information. The method provides the user an easy access to essential information and supports the expert in his task to identify the root cause of the problem.

The method is targeted, but not limited, to a daily performance analysis. The data from previous day is fed in to the system and the network elements requiring actions are shown in groups in which elements have similar problems.

The method consists of three phases, which classify the network elements by their type, performance and behaviour.

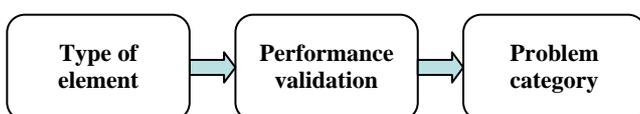


Fig. 2. Three-phase method to satisfy the requirements.

3.1. Preprocessing

Preprocessing is an essential part of data mining and analysis process. Preprocessing affects the results of the actual analysis [3, 13, 14].

Data cleaning and scaling are the preprocessing steps performed to the data before starting actual analysis. Data cleaning removes clearly erroneous values and handles missing values. Erroneous values are those outside the measurement range. Many performance indicators are on percentage scale, thus the range is from 0 to 100 and any values outside this range are errors. Erroneous values may originate from errors or malfunctions in data retrieval or storage. Erroneous values are removed from the data and treated as missing values.

The decision support system should be able to run the analysis even if the data is not complete and some values are missing. Samples with values missing can be neglected or the missing values can be estimated [15]. In our system the missing values affect the metrics of the state space. The calculation is performed with the available variables only.

The elements' quality indicators are scaled continuously and piecewise linearly to interval [0, 1] extremes corresponding to the worst and the best performance, respectively. The mapping was constructed on the basis of a priori information of network experts. Four values of the performance indicators are defined: *worst possible*, *very poor*, *satisfactory*, and *best possible*. These values are scaled to 0, 0.2, 0.9 and 1 respectively and the scaling function is created with linear interpolation. The scaling function parameters can be adjusted to different performance indicators, different networks and target performance levels. After the scaling all performance indicators are within same range and the same value refers to the same level of performance in each indicator.

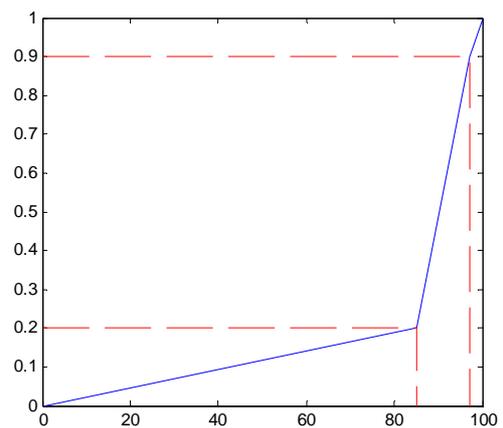


Fig. 3. Example of piecewise linear scaling function

3.2. Typification of the Network Elements

The first actual analysis phase classifies elements by their functional type within the network. Examples of types categories are indoor/outdoor cells or cells along a busy highway, elements from different network domains (radio/core network), and different network types (GSM/WCDMA).

This phase enables individual treatment of all element types. For example, the performance indicators may be different and may be gathered from different sources. Also the performance target and the scaling function parameters may be different. Network operator may set different performance criteria for a cell serving hot spot location in city centre than to a cell located in rural area.

The method divides the cells by their functional type into three groups: indoor, micro and macro. The macro cells are further divided to four categories according to both total traffic and amount of handovers (HO). The mean value of traffic and HO in available history is calculated for each cell. Then 2/3 quantile of the distribution of mean values is used to threshold the cells to groups of "high traffic", "low traffic", "high handovers" and "low handovers".

The method allows the typification based on different parameters according to operator preferences and processes in performance management.

3.3. Performance Validation

Second phase filters out the elements either performing well enough, according to operator requirements, or having no major problems. This group of elements – usually constituting the majority of all elements – require no further attention and thus the amount of data relevant in third phase is drastically reduced.

The decision if an element is performing well or requiring further attention can be based on different approaches. For example, the operator may want to select a given number or percentile of the worst performing elements to be analysed in detail. Alternatively, the selection may be based on fixed performance targets or on the performance history of elements.

All cells that have all KPIs above satisfactory, (0.9 in the scaled units) are considered to be performing well enough and require no further analysis. Euclidean distance from the ideal state (unit vector in scaled space) is calculated for the rest of the cells. Ignoring the missing variables is equivalent to assuming them to have value 1 i.e. to be in the ideal state. The cells are sorted according to the distance, largest distance meaning the most serious problems. The number of cells to analyse in the problem clustering phase is limited to maximum of 15% of the total number of cells. Cells with largest distances from the ideal are forwarded to the next phase.

3.4. Problem Clustering

In the third phase, the problematic elements are clustered based on their behaviour expressed with the scaled indicators. Clustering is done by agglomerative hierarchical method with Ward linkage [16]. The optimal number of problem clusters is tested with Davies-Bouldin index [17] and the mean silhouette method [18]. Thus, the amount of problem clusters is dynamic based on the input data and the problem data in it.

The clustering reveals the typical problem categories. The problem categorisation helps the expert to continue the analysis and to identify the root causes for each problem class. The elements in problem category have similar behaviour according to the input data. If the indicators in

input data are selected properly i.e. they are such that they have high relevancy, the further actions performed by the expert can be the same for all elements in the problem cluster.

The expert can, for example, query more detailed information about the problem, generate a report and forward the problem, or change the configuration parameters of the elements.

4. CASE STUDY AND RESULTS

4.1. Data for the Case Study

The data for the case study was a network performance database from a commercial European GSM operator. The database contained the network radio performance counters, of which the most important were used. The performance data was collected from 2385 GSM radio cells. The measurement period was 6 weeks. The data was aggregated so that for each counter one sample per day was available.

The data contained numerous invalid values and missing samples. Also possible changes in network configuration is invalidating the data and causing several anomalies in measurements. This is a typical situation in network management and also the analysis system should be able to tackle errors and anomalies in data set. In the data set there were in total 14 cells that contained invalid data or the data was partially missing.

4.2. Performance Indicators

For the high-level cluster analysis of the radio network performance and of its degradations, the performance indicator set should cover the key functionalities of GSM networks. The analysis method presented in this paper was designed to support any feature vector of performance indicators, but only one feature vector was used in the case study.

The data for typification of the elements consists of:

- *Cell type*; this parameter defines if the cell is macro cell, micro cell or indoor cell.
- *Handover amount per day*; this indicator describes how many cell handovers was made into or from the cell during a day.
- *Traffic amount per day*; the amount of calls in Erlang capacity unit.

The feature vector describing each element's performance to be validated and clustered is:

- *Dropped Call Rate*; gives the percentage of calls dropped during a day.
- *Handover Success*; gives the percentage of successful handovers into or from the cells during a day.
- *Congestion*; describes how many seconds of the day the element has been in a state that no new calls could be accepted due to lack of resources.
- *RXQUAL in classes 1 to 4*; this indicator gives the percentage of measured radio quality samples in the "good" quality classes.
- *Average Downlink Signal Strength*; this is the average signal strength received by the mobiles served by the cell during a day. The unit is dBm.

- *Call Setup Success Rate*; gives the percentage share of successful calls setup processes during a day.

4.3 Typification of the Case Study Data

The amount of macro, micro and indoor cells were 2168, 191 and 26, respectively. Because the amount of micro and indoor cells was rather low, more detailed typification would not bring benefit.

The macro cells were further divided to 4 types based on the amount of traffic and handovers. The number of low-traffic and low-handover cells was 1298, that of low-traffic high-handover 147, that of high-traffic and low-handover 146, and the number of high-traffic high-handover 577.

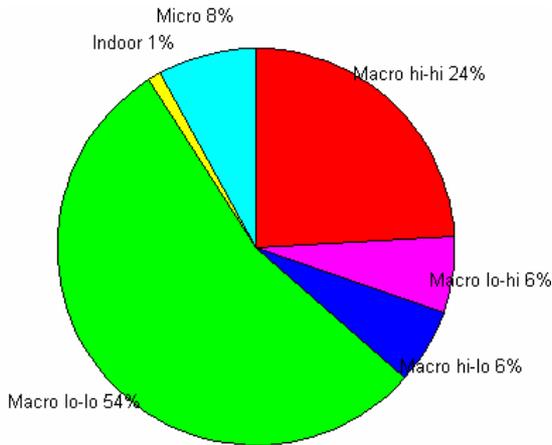


Fig. 4. The distribution of the elements after typification.

4.4 Performance Validation of the Elements

In the case study the performance verification 15% of the worst performing cells were considered as problematic and thus taken to the further analysis.

Number of problematic micro cells was 28 (14,7% of total) number of problematic indoor cells was 3 (11,5%), and number of problematic macro cells was 327 (15,1%).

4.5 Problem Clustering the Case Study Data

Table 1 shows the statistics after the third and final phase of the method. The amount of problem clusters was 5.

Table 1. The amount of cells in problem categories (I-V) after the third phase: Macro (low, low) refers to macro cell with low traffic and low amount of handovers.

	Sum	Micro	Indoor	Macro (low, low)	Macro (low, high)	Macro (high, low)	Macro (high, high)
Sum		28	3	189	18	3	117
I	37	9	1	11	1	0	15
II	9	1	0	3	2	0	3
III	21	3	0	12	1	0	5
IV	172	6	0	152	9	1	4
V	119	9	2	11	5	2	90

Fig. 5. shows the centroids of the 5 problem clusters. The graph shows that the characteristics of each problem cluster are different.

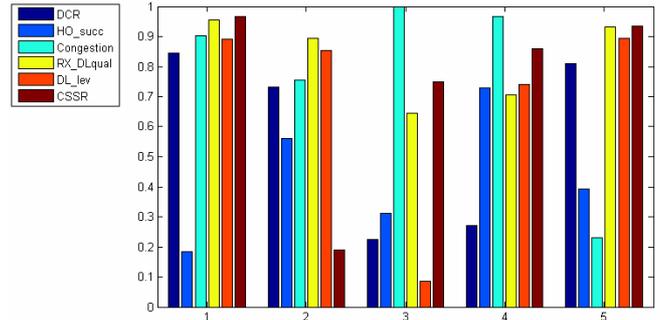


Fig. 5. The centroids of the problem clusters.

The data from previous days was associated with the clusters. History of three cells is presented in Fig. 6.

Blue asterisk (*) gives an example of a cell that has belonged to cell 0 i.e. performed well until the last day. On the last day it was assigned to problem cluster 5.

Two other cells have been members of a problem cluster most of the time.

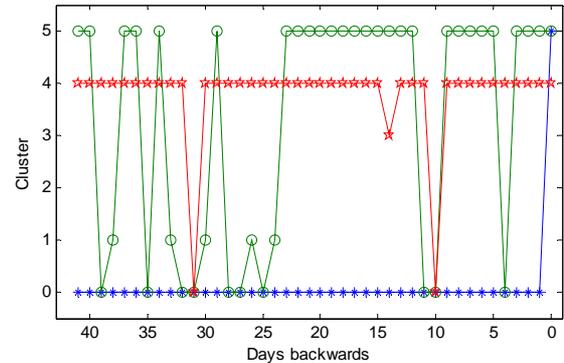


Fig. 6. Examples of problem clusters in history for three cells.

4.6 Problem Clustering Results Analysed by Expert

The boxplots [19] of the problem cluster were presented to the radio network expert. This section describes the expert's analysis of the results.

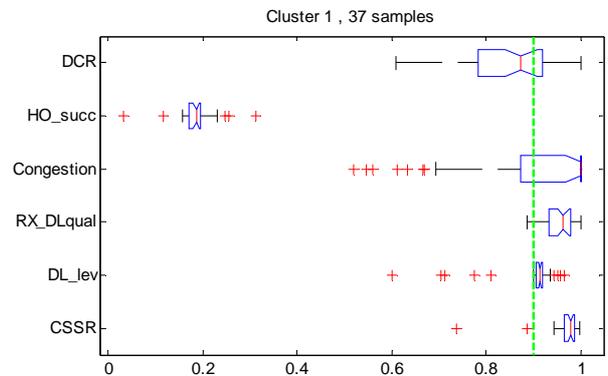


Fig. 7. Cluster 1

The problem Cluster 1 contains most probably cells with a poor adjacency plan (adjacency plan defines the neighbouring cells to which the handover is possible). As a

corrective action the definition of a revised adjacency plan should fix majority of the problems.

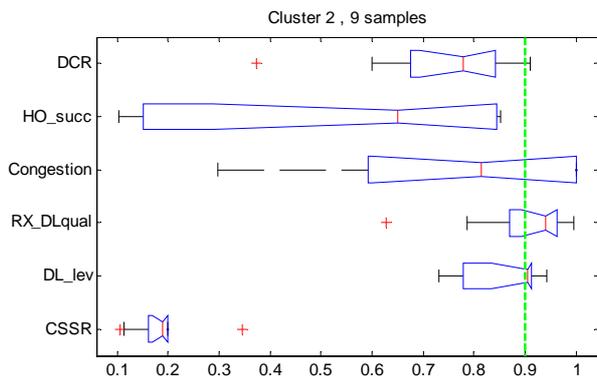


Fig. 8. Cluster 2

Because the number of cells in problem cluster 2 is low and call setup success rate is poor a possible root cause for the problems is hardware failure in the element or in the transmission links. Another possible cause is a lack of signalling capacity in the elements. The corrective action for these cells could be checking hardware and transmission links. In some cases resetting the element may solve the problems.

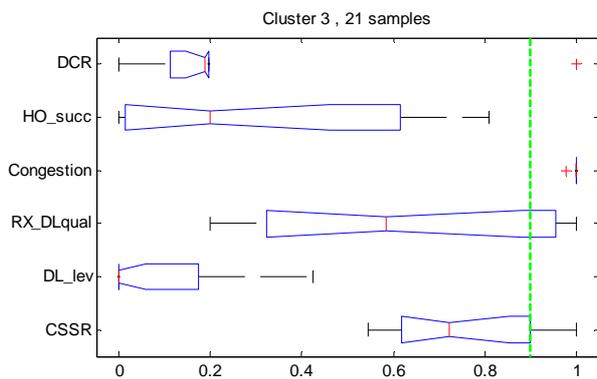


Fig. 9. Cluster 3

The problems in cluster 3 are most probably caused by poor coverage (low signal strengths). This causes some quality problems and dropped calls. To solve these problems adding new cells or changing the antenna bearings or tilts might be a solution.

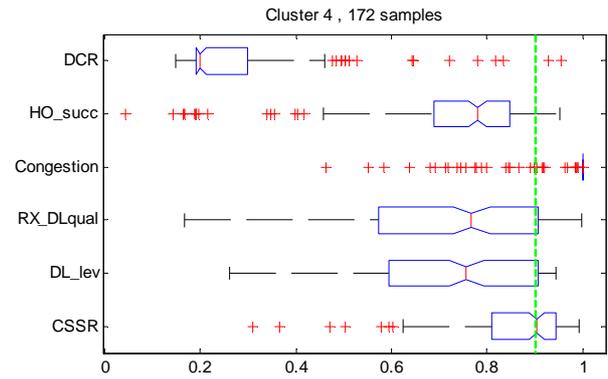


Fig. 10. Cluster 4

In cluster 4 there are slightly low signal levels, but the main root cause for the problems is poor radio quality. Typically, changing the frequencies of the problematic cells is solving these problems. In order to find the problematic frequencies, the situation should be studied on geographical map.

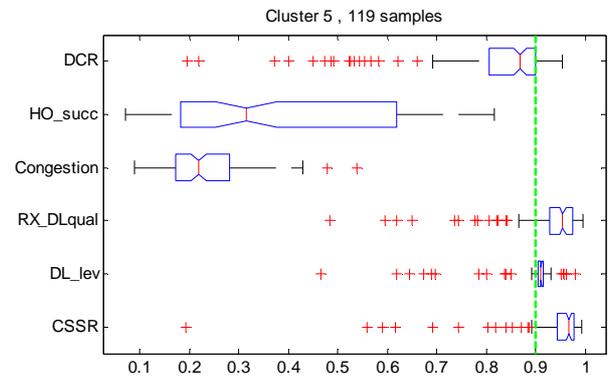


Fig. 11. Cluster 5

The problem cluster 5 is very similar to cluster 2. The difference is that call setup success rate is acceptable in cluster 5. The assumption is that the main problem in this cluster is caused by congestion, which leads to failures in handovers. Corrective action would be adding more radio capacity to elements or possibly by changing capacity configuration parameters.

5. CONCLUSIONS

The three-phase method developed in this research can be used to summarize information about a cellular mobile network. It finds well-performing elements and clusters elements with unsatisfactory performance according to the type of problem.

Experts' decision-making is supported by showing the main characteristics of problem clusters. A decision on how to improve network performance can be made collectively to several elements.

This method supports analyzing different element types with any set of indicators. In the current implementation the corner points of the scaling function for each indicator are

the only input required from the user. Those supplied, the procedure can be run fully automatically. When once defined for the network the corner points can be re-used.

Further studies are targeted towards fully automatic scaling and performance on other networks.

REFERENCES

- [1] P. Vehviläinen, K. Hätönen, P. Kumpulainen, "Data mining in quality analysis of digital mobile telecommunications network", *Proceedings of XVII IMEKO World Congress*, Dubrovnik, Croatia, June 22-27, 2003.
- [2] J. Suutarinen, "Performance Measurements of GSM Base Station System", *Thesis (Lic.Tech.)*, Tampere University of Technology, 1994.
- [3] J. Han, M. Kamber, "Data Mining: Concept and Techniques", *Morgan Kaufmann*, 2001.
- [4] O. Simula, J. Hollmén, E. Alhoniemi, "Models from data: analysis of industrial processes and telecommunication systems", *In Proceedings of the National Conference on Industrial Automation*, Ancona, Italy, pp. 13-19, November 2001.
- [5] M. Federley, E. Alhoniemi, M. Laitila, M. Suojärvi, R. Ritala, "State management for process monitoring, diagnostics and optimization", *Pulp & Paper Canada*, Volume 103, No. 2, pp. 40-43, 2002.
- [6] J. Laiho, K. Raivio, P. Lehtimäki, K. Hätönen, O. Simula: "Advanced Analysis Methods for 3G Cellular Networks". *IEEE Transactions on Wireless Communications*, (In press).
- [7] K. Julisch, "Clustering Intrusion Detection Alarms to Support Root Cause Analysis", *ACM Transactions on Information and System security*, Volume 6, No. 4, pp. 443-471, November 2003.
- [8] G. Jakobson, M. Weissman, "Real-time telecommunication network management; Extending event correlation with temporal constraints", *Integrated Network Management IV* (Chapman & Hall, London, 1995), pp. 290-301.
- [9] K. Hätönen, M. Klemettinen, H. Mannila, P. Ronkainen, H. Toivonen, "Knowledge Discovery from Telecommunication Network Alarm Databases", in: *Proc. of the 12th Int. Cong on Data Engineering*, (New Orleans, Feb. 1996).
- [10] V. Ganti, R. Ramakrishnan, J. Gehrke, A. Powell, J. French, "Clustering Large Datasets in Arbitrary Metric Spaces", *Proceedings., 15th International Conference on Data Engineering*, 1999., pp. 502 - 511, 1999.
- [11] R. Agrawal, J. Gehrke, D. Gunopulos, P. Raghavan, "Automatic Subspace Clustering of High Dimensional Data for Data Mining Applications", *Proceedings. ACM SIGMOD Int'l Conf. Management of Data*, pp. 94-105, 1998.
- [12] D.H. Fisher, "Knowledge acquisition via incremental conceptual clustering", *Machine Learning*, 2(2), 1987.
- [13] K. Hätönen, P. Kumpulainen, P. Vehviläinen, "Pre- and Post-processing for Mobile Network Performance Data", In: R. Tuokko(ed), *Automaatio03, Seminaaripäivät, Automation Makes it Work - Automaation sovellukset ja käyttökokemukset, Finnish Society of Automation*, Helsinki, pp. 311-316, September 9-11, 2003.
- [14] K. Hätönen, S. Laine, T. Similä, "Using the LogSig-function to integrate expert knowledge to Self-Organising Map (SOM) based analysis", *IEEE International Workshop on Soft Computing in Industrial Applications*, Birmingham University, New York, June 23-25, 2003.
- [15] J.L. Schafer, "Analysis of Incomplete Multivariate Data", *Chapman & Hall*, London, 1997.
- [16] R.A. Johnson, D.W. Wichern, "Applied multivariate statistical analysis 4th Edition", *Prentice Hall*, 1998.
- [17] D.L. Davies and D.W. Bouldin, "A cluster separation measure", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 1, no. 2, pp. 224-227, April 1979.
- [18] P.J. Rousseeuw, "Silhouettes: A graphical aid to the interpretation and validation of cluster analysis", *Journal of Computational and Applied Mathematics*, Volume 20, November 1987, Pages 53-65.
- [19] R. McGill, J.W. Tukey, and W.A. Larsen, "Variations of Boxplots," *The American Statistician*, Vol. 32, pp.12-16, 1978.
- [20] C. Bounsaythip, E. Rinta-Runsala, "Overview of Data Mining for Customer Behavior Modeling", *Research Report, VTT Information Technology*, 29. June 2001

AUTHOR(S): Mikko Kylväjä, Operations Solutions Business Unit, Nokia Networks, P.O. Box 370, FIN-0045 Nokia Group, Finland. Tel. +358-40-288-7105, E-mail: mikko.kylvaja@nokia.com.
Pekka Kumpulainen, Measurement and Information Technology, Tampere University of Technology, P.O. Box 692, FIN-33101, Tampere, Finland. Tel. +388-3-365-2458, fax +358-3-2171, E-mail: pekka.kumpulainen@tut.fi.
Kimmo Hätönen, Software & Applications Technologies Laboratory, Nokia Research Center, P.O. Box 407, FIN-00045 Nokia Group, Finland. Tel. +358-50-483-7278, E-mail: kimmo.hatonen@nokia.com.