

TOWARD A DISTRIBUTED DATA MINING SYSTEM FOR TOURISM INDUSTRY

Danubianu Mirela

Stefan cel Mare University of Suceava Faculty of Electrical Engineering and Computer Science 13 Universitatii Street, Suceava mdanub@eed.usv.ro mobile: +40744.547164

Socaciu Tiberiu

Stefan cel Mare University of Suceava Faculty of Economic Science and Public Administration 13 Universitatii Street, Suceava socaciu@seap.usv.ro

Bărâlă Adina

Stefan cel Mare University of Suceava Faculty of Electrical Engineering and Computer Science 13 Universitatii Street, Suceava adina@eed.usv.ro

Romania has a huge tourist's potential, but currently it is too little valued and exploited. As a result, one of the strategic developments of the economy aimed the tourism industry. The strategic decisions are based on different trends obtained from sophisticated analysis of data. Data mining helps businesses sift through layers of seemingly unrelated data for meaningful relationships, where they can anticipate, rather than simply react to, environment challenges. The aim of this paper is to present two models of data mining systems, considering that the data is processed from a distributed database.

Keywords: tourism industry, data mining techniques, distributed databases

JEL code : C89, D89

Introduction

Tourism is a great income generator due to an increased demand for its services. Their components range from quality and wide range of transportation to infrastructure, accommodation, food and beverage, support services and travel distribution services.

Romania has a huge tourist's potential, unfortunately, too little valued and exploited. As a result, one of the strategic developments of the economy aimed the tourism industry. But strategies are based on different trends obtained from sophisticated analysis of data. Providing the managers in the tourism industry with information about and insight into the existing data is the key function of the data warehouse systems [1].

Data mining - techniques for exploration and analysis of large quantities of data in order to discover meaningful patterns and rules - helps businesses sift through layers of seemingly unrelated data for meaningful relationships, where they can anticipate, rather than simply react to, environment challenges.

A system which enables the use of data mining techniques on data stored in a data warehouse is ideal for high quality analyzes to support strategic decision. But designing and implementing a Data Warehouse is a complex and expansive process [2], so we can apply the data mining algorithms on large volumes of data from relational databases.

The aim of this paper is to present the opportunity to use data mining methods on data from tourism and also to present two models of data mining systems, considering that the data is processed from a distributed database.

Data mining

Data mining is the exploration and analysis, by automatic or semiautomatic means, of large quantities of data in order to discover meaningful patterns and rules. So, data mining is defined as the process of extracting interesting and previously unknown information from data, and it is widely accepted to be a single phase in a complex process known as Knowledge Discovery in Databases (KDD).

This process consists of a sequence of the following steps [3]:

data cleaning – to remove noise and irrelevant data

data integration – where multiple data sources are combined

data selection – for retrieve from the database only the relevant data for the analyze

data transformation – where data are transformed or consolidated into forms appropriate for mining

data mining – the phase where the algorithms are applied in order to extract data patterns

pattern evaluation – to find the interesting patterns who representing new knowledge

knowledge presentation – when the visualization techniques are used to present the mined knowledge to the user

In order to ensure that the extracted information generated by the data mining algorithms is useful, additional activities are required, like incorporating appropriate prior knowledge and proper interpretation of the data mining results.

Figure 1 presents these phases and the most important interdependencies between them.

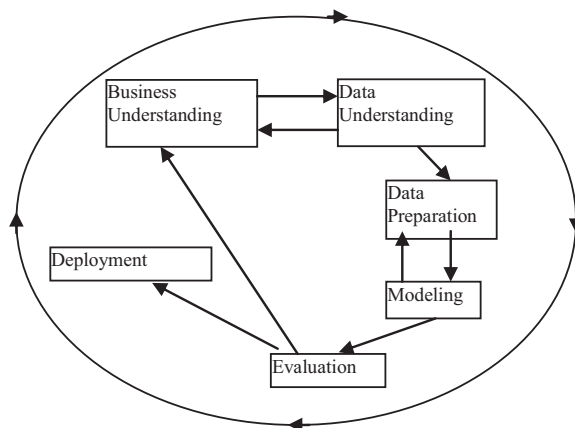


Figure 1. The steps of the CRISP-DM process (adapted from Wirth and Hipp)

Information Systems Used in Tourism Industry

Information technology was initially viewed by the tourism industry as a back-office function that supports the finance and accounting areas. The industry has advanced far beyond this view during the past decade. As information is vital for tourism industry, effective use of Information Technology is necessary.

An Information System regarding tourism activity should have some characteristics. It should collect, select and process information that is internal to the tourism activity coming from entities related to this sector such as National Tourism Agency and it should have subsystems that receive data from other business sectors, such as passengers landed from regular and nonregular flights, passengers flying with low-cost. It should support the decision-making process via integration of data from various sources, integrating them in a manner that permits analyses and comparisons between tourism indicators for the different regions of the countries and for different time periods. Also it should include historic information, such as the amounts of overnight stays per country of residence and for all types of hotel establishments throughout the 12 months of the year for every region [4]. It should also support the specialists who process and use tourism demand forecasts. As such, it detects the 10 countries with the highest tourism demand in the previous year. Taking into account the 10 countries listed and the forecasting methods suited to the specificity and nature of the data, the system prepares the necessary data to be processed by the specialists who create forecasts for tourism demand. [5]

Currently, the most used information systems in Romanian tourism industry are the *front-office systems and the reservation systems*.

Front-office information systems are those data processing systems that provide reports in visual or written form. They are used mainly in the management of tourist accommodation (hotels, motels, hostels or cruise ships) or in the travel agencies activities. These systems may be used for: tourists registration when the personal data about tourists are collected; marketing of various tourism products, such as rental cars; rooms management, when are collected and processed data regarding the rooms status, (allows instant viewing of room availability for all room types, indicates whether rooms are dirty or clean, allows rooms to be placed out of inventory or out of order to restrict rental) and tracks of revenues, providing transaction processing and obtain information about any debts and credits in relation to customers

Information Systems Used for Reservations provide rapid access to information and ensures the accuracy of this information. They bring information services, booking and selling and are used both by individual tourists and travel agents or commissioners. Most often this type of systems uses Web technologies. These systems use hardware and software specific to conduct them activities. Although providers of tourist services in Romania currently use such systems for ticketing most, is well to remember that these systems can be used for marketing or management activities.

In the tourism industry knowing the guests - where they are from, how much they spend, and when and on what they spend it- can help a company to formulate marketing strategies and maximize profits. Due to technological development touristic companies have accumulated large amounts of customer data, which can be organized and integrated in databases that can be used to guide marketing decision [6]. Since identification of important variables and relationships located in these consumer-information systems can be a difficult task, some companies have attempted to raise the power of information by using *data mining technologies*.

How can use Data Mining Technologies in Hospitality

Hospitality is used to describe hotels and similar accommodations as well as restaurants and catering organizations (Holloway, 2006) and represent a very important aspect of the tourist industry.

If hospitality organizations want to compete successfully, they must do so by using technology to drive value to both the customer and to the firm.”[]

In this area Information Systems have been used to assist the delivery of hospitality services. Some of the key ways are (Buhalis, 2003): improved capacity management and operations efficiency; central room inventory control; last room availability information; yield management capability; marketing, sales and operational reports; tracking frequency flyers and repeat hotel guests; internal management of operations from transactions to human resources.

Most of the items on the above list apply only to hotels and accommodation providers

In order to make high quality marketing research and planning data-mining technology allows hotel companies to predict consumer-behavior trends, which are potentially useful for marketing applications.

The tasks performed by data mining can be grouped into the following five categories.

Classification arranges customers into pre-defined segments that allow the size and structure of market groups to be monitored. Also, predictive models can be built to classify activities. Classification uses the information contained in sets of predictor variables, such as demographic and lifestyle data, to assign customers to segments.

Clustering group customers based on domain knowledge and the database, but does not rely on predetermined group definitions. This function is beneficial because it aids hoteliers in understanding who their customers are. For example, clustering may reveal a subgroup within a predetermined segment with homogenous purchasing behavior (a subgroup of holiday shoppers within the transient segment) that can be targeted effectively through a specific ad campaign with the scope that the members of the subgroup will increase their number of stays or become more loyal. On the other hand, clustering may indicate that previously determined segments are not parsimonious and should be consolidated to increase advertising efficiency. Information such as demographic characteristics, lifestyle descriptors, and actual product purchases are typically used in clustering.

Deviation detection uncovers data anomalies, such as a sudden increase in purchases by a customer. Information of this type can prove useful if a hotel corporation wants to thank a guest for her or his recent increase in spending or offer a promotion in appreciation. Marketing managers may also attempt to draw correlations between surges in deviations with uncontrollable business-environment factors that are not represented in the database.

Association entails the detection of connections between records, driven by association and sequence discovery. For example, a possible detected association may be that a particular segment’s average length of stay increases after a specific advertising campaign. Another association task could be employed in an effort to determine why a specific promotion was successful in one market, but ineffective elsewhere. Specific information regarding customer-purchase histories is necessary to formulate probabilistic rules pertaining to subsequent purchases.

Forecasting predicts the future value of continuous variables based on patterns and trends within the data. For instance, the forecasting function can be used to predict the future size of market segments. With forecasting one can also use data trends to project which hotel amenities are of growing importance to consumers and will be key drivers of the future perception of value of consumers.

Distributed Model for a Data Mining System

Data Mining systems have the following characteristics: they must not limit the size of data sets, the performances are optimized for large data sets and they are enough flexible to use various techniques of data mining. Also they offer support for multi-user access and requires a total control over data access. Finally they provide management and maintenance at a distance.

The basic elements of a data mining system of data are: user interface, the specific data mining services, data access and data itself.

Usually, Data Mining systems are built using client-server architecture, with different distribution on the two components of the items listed above.

In order to achieve a prototype of a data mining system for hotel industry of Bucovina, we proposed two models of system architecture and we have studied some of their advantages and disadvantages.

We started from the reality that each accommodation establishment manages its own data. Passing over specific needs all these work with databases containing data on customers, on services requested, on the amount spent, so on... If these systems allow a part of their data, in terms of ensuring data privacy, to be used for analysis, then it is possible that projections by necessary attributes of the tables to be available for sharing. If individual systems are connected through a communications network can assume that we are dealing with a heterogeneous distributed database, as shown in Figure 2.

In order to study the two models we chose ten accommodation establishments that have agreed to share the data available for this purpose.

The first stage was a selection of required data. For that we have applied a projection by a list of fields with the same meanings of the tables, corresponding to the following model:

$$\prod_{c_{1k}, c_{2k}, \dots, c_{nk}} (T_k) \quad (1)$$

Where $(c_{1k}, c_{2k}, \dots, c_{nk})$ represent the list of n attributes required from table T_k .

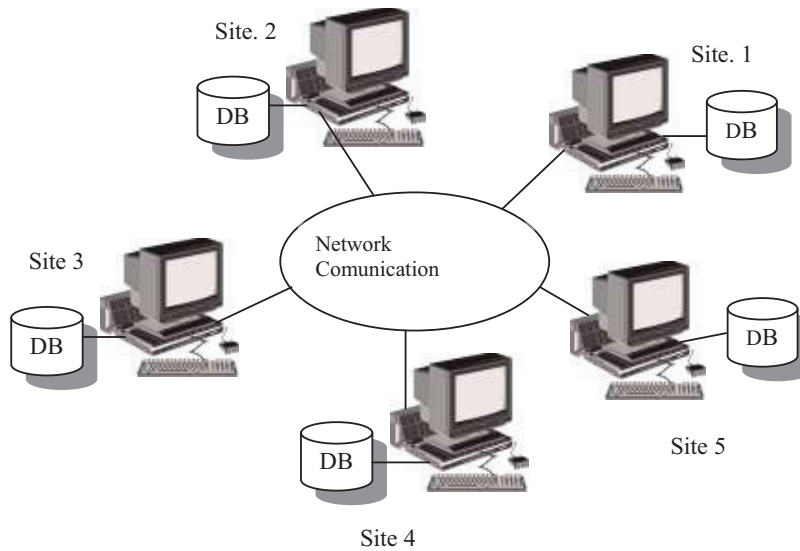


Figure 2. Distributed Database Architecture

The set of all these projections in a site form a fragment of the distributed database. Starting from this condition were analyzed the following situations.

In the first case, were installed in each site user interfaces and services suitable for data mining. Thus it was possible to apply local data mining methods (e.g. discovery of association rules). Local results were replicated in a single node and were combined to obtain the global solution. It is obvious that in the site where data were replicated were required additional operations in order to verify and validate the final results.

For example, for association rules was necessary to calculate the global support and confidence. In this case there are two major disadvantages. The first is related to the small volume of data processed on the local sites, which may lead to partial results inconclusive. The second disadvantage is the need to conduct further operations in the site where the results are collected.

The other option is a replication of all fragments in a single site where the specific components for data mining systems are installed. On these aggregated data we apply different methods of data mining.

The advantage of this approach lies in that additional operations for further validation of the global results are eliminated. However, there is a drawback related to large volume of data transferred on the network.

Conclusion and Future Work

In this paper we have shown that data mining techniques can be applied successfully in the field of tourism, especially in connection with strategic marketing. We also kept in mind that these techniques can be applied to data sets from various sources, which can be successfully treated as fragments of a distributed database. In this context we have examined two models of distribution of components of a data mining system and how we can apply specific methods and we have underline their advantages and disadvantages.

References

1. M. Danubianu, T. Socaciu, A. Barila Some Aspects Of Data Warehousing In Tourism Industry, accepted for publishing in The Annals of the "Stefan cel Mare" University Suceava. Fascicle of The Faculty of Economics and Public Administration, 2009
2. M. Danubianu, Advanced Information Technology – Support of Strategic Decision in Romanian Tourism Industry , in Proceedings of IECS 2009, Sibiu , Romania
3. R. Wirth and Hipp, (2000) J. CRISP-DM: Towards a standard process model for data mining. In Proceedings of the 4th International Conference on the Practical Applications of Knowledge Discovery and Data Mining, pages 29-39, Manchester, UK
4. A. Poon Tourism, Technology and Competitive Strategies. Wallingford: CAB International, 1993
5. C. Ramos, F. Perna, Information system for Tourism Activity Monitoring and Forecasting Indicators as an experience for Portugal, Tourism and Hospitality Research, February, 2009

6. M. Danubianu, V. Hapenciu Improving Customer Relationship Management In hotel industry by Data Mining Techniques, Proceeding of “Competitiveness and Stability in the Knowledge-Based Economy”, Vol: CD, 30-31 Mai, 2008, Craiova, Romania, ISSN/ISBN: 978-606-510-162-3, Pagini: 2444-2452