

Sequence Data Mining for Adverse Event Prediction and Action Recommendation

Dr. Reza Sanati-Mehrizy, Utah Valley University

Reza Sanati-Mehrizy is a professor of Computer Science Department at Utah Valley University, Orem, Utah. He received his M.S. and Ph.D. in Computer Science from the University of Oklahoma, Norman, Oklahoma. His research focuses on diverse areas such as: Database Design, Data Structures, Artificial Intelligence, Robotics, Computer Aided Manufacturing, Data Mining, Data Warehousing, and Machine Learning.

Dr. Afsaneh Minaie, Utah Valley University

Afsaneh Minaie is a professor of Computer Engineering at Utah Valley University. She received her B.S., M.S., and Ph.D. all in Electrical Engineering from University of Oklahoma. Her research interests include gender issues in the academic sciences and engineering fields, Embedded Systems Design, Mobile Computing, Wireless Sensor Networks, and Databases.

Dr. Ruhul H Kuddus, Utah Valley University

I obtained my Undergraduate degree from University of Dhaka, Dhaka, Bangladesh; MS in Biology from George Mason University, Fairfax VA; and Ph.D. in Molecular Genetics and Biochemistry from University of Pittsburgh, Pittsburgh, PA. My research area include biomarkers in molecular medicine, cancer epidemiology and organ transplantation. Recently I also included effects of climate change on public health in my research agenda. My research also involve data mining.

Dr. Ali Sanati-Mehrizy

Dr. Ali Sanati-Mehrizy is a Pediatric resident physician at Rutgers University - New Jersey Medical School in Newark, NJ. He is a graduate of the Milton S. Hershey Pennsylvania State University College of Medicine. He completed his undergraduate studies in Biology from the University of Utah. His research interests are varied and involve pediatric hematology and oncology as well as higher education curricula, both with universities and medical schools.

Mr. Paymon Sanati-Mehrizy, Icahn School of Medicine at Mount Sinai

Paymon is currently a medical student at the Icahn School of Medicine at Mount Sinai. He completed his undergraduate studies in Biology from the University of Pennsylvania in May 2012. Currently, his research interests consist of higher education curricula, both with universities and medical schools.

Sequence Data Mining for Adverse Event Detection and Action Recommendation

Abstract

Many real-life data mining applications use sequence data modeling in which data is represented as a sequence. A temporal sequence is finite ordered list of events $(t_1, e_1), (t_2, e_2), \dots, (t_n, e_n)$ where t_i represents time and e_i represents the event taking place at time t_i . e_i takes place before e_{i+1} for $1 \leq i \leq n-1$. This model can be used in data mining, called sequence data mining, to predict certain event that may take place at a specific time. Sequence data mining has a wide range of applications. This data mining technique can be used for prediction of adverse events and recommend proper actions to be taken as needed. In aviation safety, the future of a flight can be predicted as a sequence and proper action can be recommended to avoid dangerous situations that a flight may get into otherwise. In health care system, the future of a bacterial infection can be predicted and proper medicine can be prescribed for different situations to bring the patient's illness to an end. In the marketing, customer shopping can be monitored and certain action can be taken, such as mailing coupons, to encourage the customer for further shopping of relevant products. In the real-life situations such as manufacturing, sensors' data can be analyzed to control operations and predict dangerous situations and recommend proper actions to be taken.

This paper discusses sequence representation, implementation, and its application for a number of different cases.

Introduction

Data mining combines tools from statistics and artificial intelligence (such as neural networks and machine learning) with database management to analyze large digital collections, known as data sets¹. It is a well-researched area of computer science with high demand because it is useful in any field where there large quantities of data from which to extract meaningful patterns and rules². Therefore, many organizations and businesses can benefit from data mining techniques as they record lots of data daily.

Sequence data mining is a specific area of data mining that has a wide range of application in this field^{3,4}. In some application, sequence data mining can be used to identify anomalous events and recommend proper actions to be taken to avoid such situations. In aviation safety, for example, the future of a flight can be predicted as a sequence of events. If a sequence of events appears anomalous, proper action(s) can be recommended to take place to avoid dangerous situations that this flight may otherwise get into. In health care system, physicians can predict the future of a bacterial infection or an allergic. It is desirable to bring the patient's illness to an end as soon as possible. To do so, the physicians prescribe proper medication considering the situations. In other real-life situations such as operating manufacturing plants, sensors' data can be analyzed to control operations, predict dangerous situations and recommend and implement proper

actions. This type of data analysis can help engineers to estimate the remaining life of equipment and recommend proper maintenance services before the equipment malfunction and the entire production line shuts down.

Definition

A sequence is a nonempty ordered list of tuples $(t_1, e_1, a_1) \dots (t_n, e_n, a_n)$ where t_i represents a time point, e_i represents an event and a_i represents an action at that time point. If the event e_{i+1} exists, it is the effect of event e_i . At any time point t_i , the event e_i may cause a set of possibly empty event(s). Any of these new events can initiate another sequence. A sequence can come to end if no new event is generated. In some applications, it is beneficial to bring a sequence to a halt situation. But in some other cases, it is desirable that an event causes another event that may initiate another sequence. At any time t_i a sequence may be brought to a halt situation if a proper action a_i is taken or it may come to an end without the need for an action taken by an outsider. In this case a_i is null.

Classification of Sequences

Han et al.⁵ have given a good coverage of sequence mining including classification of sequences. The classification depends on what criteria are used for classification. In their classification, sequences are considered as complex data type.

In this paper, the sequences are simply classified as desirable sequences and undesirable sequences just for the purpose of this research work. Desirable sequences are those that their existence is beneficial and we want them to exist. Undesirable sequences are those that their existence is harmful and they need to be brought to an end as soon as possible. In marketing, customers' shopping behavior can be monitored and certain action can be taken, such as mailing coupons, to encourage the customer to continue shopping of relevant products that actually may initiate new sequence(s). A radio station may follow the listening habits of clients of different age groups and reward the groups by adjusting the broadcasting line-ups and infuse appropriate advertisements in order to boost profit of the corporations. In these cases, the sequence is enhanced and augmented for expected positive outcomes. In health care system, physicians can predict the future of a bacterial infection or an allergic. These types of sequences are harmful and need to be brought to an end as soon as possible. To do so, physicians recommend necessary treatments to terminate these undesirable sequences.

Implementation

A multiway lexicographic search tree can be used to represent a sequence where an event from the sequence of events determines a multiway branch at each step. If the sequence is constructed from the English alphabets, at the root of the tree there are 26 possible branches followed by another branch according to the next letter in the sequence.

The tree shown in Fig. 1⁶ represents sequences constructed from letters a, b, and c.

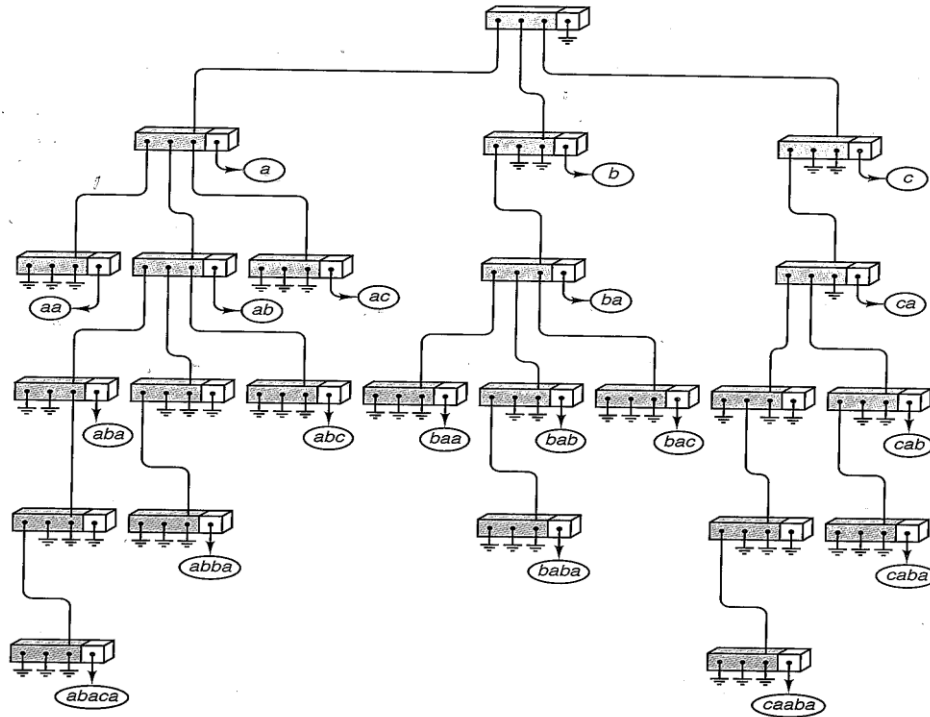


Fig. 1. A multiway lexicographic search tree constructed using the English letters.

At the root of the tree, there are three branches one for each of the letters a, b, and c. At the next level there are three nodes and each node represents three branches similarly. For information retrieval, where these sequences represent keys of records, each node contains a pointer to the record of information with this sequence as its key. This data structure can be used to represent the sequences of events where each branch of a node represents an event and each path in the tree represents a sequence of events. When a path terminates, it means that the corresponding sequence of events came to an end. Each node of the tree can be considered the outcome of a sequence at a particular time. The concatenation of events along the path from the root of the tree to this node is the cause of this outcome. At any node of the tree, the outcome can be evaluated. If the outcome is undesirable, we can look back for the cause of it and tune the system for a better outcome.

Application

For the purpose of sequence mining where a sequence is constructed from events, nodes contain more branches and each node contains a pointer to an action (possibly a set of actions) that needs to be taken place for the event leading to that node. This action may prevent initiation of any further event or may promote generation of new sequences. In the following diagram shown in Fig. 2, a customer has purchased baby food, coupon for baby diaper can be sent to him that more likely encourages him to purchase baby diaper. Later on when baby gets older, coupon for toys and bicycle can be sent to him and he again may buy toys and bicycle. Each of these events is a starting point of a sequence.

Having such sequence presented, when a customer purchase baby food, we can predict that he will purchase baby diaper, toys and bicycle at some point in time.



Fig. 2. The sequence of events and action recommendations in a customer shopping monitoring system.

In terms of health care system, sequence data mining can be useful in patient management to operating a hospital network. Let us consider a patient showing signs of type I hypersensitivity (allergic) reaction due to inhalation of an allergen. The person is having frequent sneezes, soon the person may develop icy throat, followed by runny nose and watery eyes, followed by asthmatic conditions followed by hay fever. If left untreated for more than a day, the patient may suffer delayed reactions of allergy that involves inflammatory reactions leading to itchy throat and skin, damage in the trachea and the alveoli of the lung. In a few days the person may get secondary bacterial infections of the lung that may lead to bronchiolitis, followed by high fever followed by pneumonitis and pneumonia. If left untreated, the person may develop septicemia, followed by vascular shock, followed by coma and ultimately death. This is an ordered sequence where each of these symptoms appears after certain period of time. This is an undesirable sequence and needs to be brought to a halt situation as soon as possible. For example interrupting the sequence by treating the subject with epinephrine when the patient developed runny nose and watery eyes could have stopped the sequence⁷. Had the sequence progressed beyond that step, specific actions could have been taken at the subsequent steps. For example, the patients could be treated with anti-inflammatory drugs within a day of the beginning of the episode or treated with antibiotic if secondary bacterial infection is diagnosed. If the recommended action is effective, it will cure the patient and terminates the sequence. If not, the progress continues and goes to the next step. In this situation the prediction is easier because we can tell what the next stem can be. Suppose the allergen was a virus and sometimes the virus mutates and changes its course and initiates a new sequence and will be harder to control the illness because the predictions become harder.

In some cases, we recommend an action to change the sequence intentionally. Consider a case when a student fails the first quiz and misses submitting the first project. The instructor can predict that this student will more likely have problem with other assignments, quizzes, exams and so on. The instructor may even guess that this student will be one of those who may fail the class. This is an example of undesirable sequence.

The instructor may take preemptive actions such as contact the student and tries to find out what the problem might be and makes appropriate recommendations. If the recommendation is effective and the students changes his behavior, the sequence changes and hopefully the student may enter in more desirable sequence leading to successful completion of the course.

Adverse Event Detection

Based on previous sequences, a model can be designed for detecting future anomalous sequences. Nikunj C. Oza⁸ presented a paper that includes a simple model for anomalous sequence detection. Consider the following sequences:

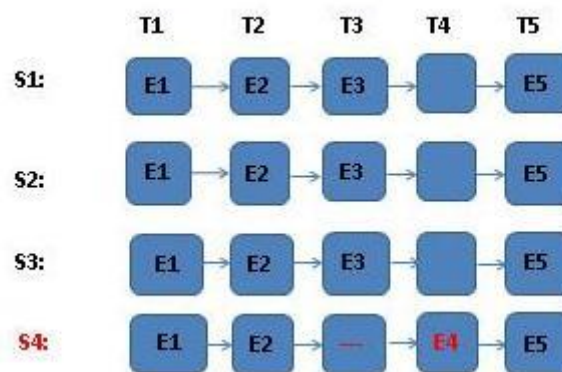


Fig. 3. Sequences

S1, S2, and S3 are normal sequences. But S4 is an anomalous sequence. Because it is expected that the event E3 take place at time t3 which did not happen and event E4 took place at time t4 which was not expected. This model can be used to detect anomalous flights in aviation safety system where events can be on/off positions of various switches.

Teaching Sequence Mining

In our computer science department, we teach an undergraduate data mining course. The major parts of workload for this class are exercise problems, data analysis projects, and research projects. Exercise problems and data analysis projects are individual works. For data analysis projects the students use Weka⁹, a popular open source suite of machine learning software for data mining. But the research projects are team works usually with two team members on each team. Each team is given certain topic in data mining to do research on. Each team has to present the result of its research to the class at the end of semester. This way, the teams share their work with others and learn from each other. Also, the students learn how to do a presentation in front of a good crud and answer questions. Naturally, this type of workload prepares student better for this highly competitive job market.

Finally, the teams will be encouraged to convert the result of their research project to a paper and submit it to a conference for presentation and publication. One little problem

that I feel need to mention is the travel expense for the student to the conference. To solve this problem, usually the students can sign up for a few hours of volunteer work at the conference and get their registration fee waved. Sometimes, the students are willing to pay for a big part of the travel expenses because this trip becomes a vacation for them and also get a publication for their resume.

In order to enrich the content of this data mining course, the plan is that two weeks of teaching sequence mining to be added to the content of this course. This gives the student a considerable experience with mining such a complex data type that has applications in almost any field.

Conclusion

Sequence data mining has multitude of potential applications in diverse disciplines, from aviation safety to health care and from student management to consumer behavior management. Careful observation and deductive reasoning may provide the basis of developing an algorithm to develop software that can be used in interrupting undesirable sequence and enhance and augment desirable sequence. In this paper, a few areas of application were presented very briefly. Finally, a simple example was included that gives an idea for detecting anomalous sequences.

References

1. <http://www.britannica.com/EBchecked/topic/1056150/data-mining>.
2. Isinkaye O. Flasadé, Computational Intelligence in Data Mining and Prospect in Telecommunication Industry, *Journal of Emerging Trends in Engineering and Applied Science*, (JETES) 2 (4): pp. 601-605.
3. Mabroukeh NR, Ezeife CI. A taxonomy of sequential pattern mining algorithms. *ACM Computing Surveys* 2010; 43:1.
4. Han J, Cheng H, Xin D, Yan X., Frequent pattern mining: current status and future directions. *Data Mining and Knowledge Discovery* 2007; 15 (1): 55–86.
5. Han et. al., *Data Mining Concept and Technology*, Third Editio, Morgan Kaufmann, 2012.
6. Robert, L. Kruse and Alexander J. Ryba, *Data Structures and Program Design in C++*, Prentice Hall, 1999.
7. Dugdale DC, Henochowicz SI, Zieve D. Allergic reactions. *ADAM Medical Encyclopedia* (2013). <http://www.ncbi.nlm.nih.gov/pubmedhealth/PMH0001076/> (retrieved Jan 2014).
8. Nikunj, C. Oza, *Data Mining for Aviation Safety*, American Statistical Association, San Francisco Bay Area Chapter, October 21, 2010.

9. [http://en.wikipedia.org/wiki/Weka_\(machine_learning\)](http://en.wikipedia.org/wiki/Weka_(machine_learning))