



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 5, Issue 2, February 2017

Securely Mining User Sequential Pattern Matching in Document Streams

K. Radhika¹, S J Sowjanya², Yerragudipadu Subbarayudu³

Assistant Professor, Department of CSE, Institute of Aeronautical Engineering Hyderabad, India¹

Assistant Professor, Department of CSE, Institute of Aeronautical Engineering Hyderabad, India²

Assistant Professor, Department of CSE, Institute of Aeronautical Engineering Hyderabad, India³

ABSTRACT: The process of data mining produces various patterns from a given data source. The most recognized data mining tasks are the process of discovering frequent item sets, frequent sequential patterns, frequent sequential rules and frequent association rules. Numerous efficient algorithms have been proposed to do the above processes. In my research work Textual documents created and distributed on the Internet are ever changing in various forms in order to characterize and detect personalized and abnormal behaviors of Internet users, proposing Sequential Topic Patterns (STPs) and formulate the problem of mining User-aware Rare Sequential Topic Patterns (URSTPs) in document streams on the Internet. They are rare on the whole but relatively frequent for specific users, so can be applied in many real-life scenarios, such as real-time monitoring on abnormal user behaviors. Research work consider three phases: preprocessing to extract probabilistic topics and identify sessions for different users, generating all the STP candidates with (expected) support values for each user by pattern-growth, and selecting URSTPs by making useraware rarity analysis on derived STPs. Experiments shows that our approach can indeed discover special users and interpretable URSTPs effectively.

KEYWORDS: KDD, User-aware Rare STPs Cross Layer, Failure Recovery, Link Stability

I. INTRODUCTION

The main goal of mining is to discover the useful and unexpected patterns in database. There is lot of work in the field of data mining about pattern mining. In this, we will be interested by a specific type of database called sequence databases. Sequence database contains some sequences. Here we find out sequential pattern in document stream. Sequential pattern mining has wide applications on the client purchase behaviour analysis, web-log analysis and medical record analysis. We find out the pattern that is frequently used by the user. This technique is useful to find out the users abnormal behaviour on the internet. Sequence database sequence pattern mining is the task of finding patterns which are present in a certain number of instances of data. The identified patterns are expressed in terms of sub sequences of the data sequences and expressed in an order that is the order of the elements of the pattern should be respected in all instances where it appears. If the pattern is considered to be frequent if it appears in a number of instances above a given threshold value, usually defined by the user, then it is considered to be frequent. There may be huge number of possible sequential patterns in a large database. Sequential pattern mining identifies whether any relationship occurs in between the sequential events. The sequential patterns that occur in particular individual items can be found and also the sequential patterns between different items can be found. The number of sequences can be very large, and also the users have different interests and requirements. If the most interesting sequential patterns are to be obtained, usually a minimum support is pre-defined by the users. By using the minimum support, sequential patterns which are not so important is taken out and hence the mining process will be more efficient Textual documents made and appropriated on the Internet are perpetually changing in different structures. The vast majority of existing works are committed to theme demonstrating and the development of individual topics, while consecutive relations of topics in progressive documents distributed by a particular user are overlooked. In this paper, so as to describe and recognize customized and abnormal behaviors of Internet users, we propose Sequential Topic Patterns (STPs) and plan the issue of mining Usermindful Rare Sequential Topic Patterns (URSTPs) in document streams on the Internet. They are



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 5, Issue 2, February 2017

uncommon overall yet generally visit for particular users, so can be connected in some real-world situations, for example, realtime observing on abnormal user behaviors. We exhibit a gathering of algorithms to explain this imaginative mining issue through three stages: preprocessing to separate probabilistic topics and distinguish sessions for various users, creating all the STP applicants with (expected) bolster values for every user by example development, and selecting URSTPs by making user-aware rarity analysis on determined STPs Document streams are made and appropriated in different frames on the Internet, for example, news streams, messages, small scale blog articles, talking messages, research paper chronicles, web gathering examinations, et cetera. The content of these documents for the most part focus on a few particular topics, which reflect disconnected get-togethers and users' qualities, all things considered. To mine these bits of data, a great deal of scrutinizes of content mining concentrated on extricating topics from archive accumulations and document streams through different probabilistic theme models, for example, established PLSI, LDA and their augmentations. Exploiting these separated topics in document streams, the vast majority of existing works broke down the development of individual topics to distinguish and predict social events as well as user behaviors. Notwithstanding, few investigates paid consideration on the relationships among various topics showing up in progressive archives distributed by a particular user, so some covered up however huge data to uncover personalized behaviors has been ignored. Knowledge discovery is a process of nontrivial extraction of information from large databases, information that is unknown and useful for user. Data mining is the first and essential step in the process of knowledge discovery. Various data mining methods are available such as association rule mining, sequential pattern mining, closed pattern mining and frequent item set mining to perform different knowledge discovery tasks. Effective use of discovered patterns is a research issue. Proposed system is implemented using different data mining methods for knowledge discovery. Text mining is a method of retrieving useful information from a large amount of digital text data. It is therefore crucial that a good text mining model should retrieve the information according to the user requirement. Traditional Information Retrieval (IR) has same objective of automatically retrieving as many relevant documents as possible, whilst filtering out irrelevant documents at the same time. However, IR-based systems do not provide users with what they really need. Many text mining methods have been developed for retrieving useful information for users. Most text mining methods use keyword based approaches, whereas others choose the phrase method to construct a text representation for a set of documents. The phrase-based approaches perform better than the keyword-based as it is considered that more information is carried by a phrase than by a single term. New studies have been focusing on finding better text representatives from a textual data collection. One solution is to use data mining methods, such as sequential pattern mining for Text mining. Such data mining-based methods use concepts of closed sequential patterns and non-closed patterns to decrease the feature set size by removing noisy patterns. New method, Pattern Discovery Model for the purpose of effectively using discovered patterns is proposed. Proposed system is evaluated the measures of patterns using pattern deploying process as well as finds patterns from the negative training examples using pattern Evolving process

II. RELATED WORK

1. Discovery of rare sequential topic patterns in document stream[1]

From This Paper I Referred Plain text documents made and circulated on the Internet are constantly changing in different structures. Mining topics of these archives has huge applications in numerous areas. A large portion of the writing is committed to point displaying, while successive examples of topics in archive streams are disregarded. Also, conventional consecutive example mining calculations basically centered around successive examples for deterministic information sets, and in this way not appropriate for document streams with topic uncertainty and uncommon examples. In this paper [1], author figure and handle the mining issue of uncommon Sequential Topic Patterns (STPs) for Internet document streams, which are uncommon all in all yet moderately regularly for particular clients, so likewise intriguing. Since this kind of uncommon STPs mirrors clients' particular practices, our work can be connected in numerous fields, for example, customized setting mindful proposal and ongoing checking on irregular client practices on the Internet. author propose a novel way to deal with finding client related uncommon STPs in light of the fleeting and probabilistic data of concerned topics. Subsequent to extricating topics from archives by LDA and sorting the record stream into sessions for various clients amid various eras, the proposed calculations find uncommon STPs by



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 5, Issue 2, February 2017

(1) digging STP possibility for every client through a proficient calculation in view of example development, and (2) creating client related uncommon STPs by example irregularity examination.

2. Mining probabilistically frequent sequential patterns in large uncertain databases[2] From This Paper I Referred Information uncertainty is characteristic in some real - world applications, for example, natural observation and versatile following. Mining successive examples from wrong information, for example, those information emerging from sensor readings and GPS directions, is vital for finding concealed learning in such applications. In this paper, author proposes to gauge design recurrence in view of the conceivable world semantics. Author[2] build up two dubious grouping information models dreamy from some real - world applications including indeterminate succession information, and figure the issue of mining probabilistically visit consecutive examples (or p - FSPs) from information that adjust to our models. Be that as it may, the quantity of conceivable universes is amazingly substantial, which makes the mining restrictively costly. Propelled by the well-known Prefix Span calculation, Author create two new calculations, on the whole called U - PrefixSpan, for p - FSP mining. U - PrefixSpan successfully stays away from the issue of "conceivable universes blast", and when joined with our four pruning and approving techniques, accomplishes shockingly better execution. Author additionally proposes a quick approving strategy to further accelerate our U - Prefix Span calculation. The proficiency and adequacy of U - PrefixSpan are checked through broad investigations on both real - world and engineered datasets. 3. Mining probabilistic frequent spatio - temporal sequential patterns with gap constraints from uncertain database

[3] From This Paper I Referred Uncertainty is regular in real - world applications, for instance, in sensor organizes and moving article following , bringing about much enthusiasm for thing set digging for questionable exchange databases. In this paper, author concentrate on example digging for dubious groupings and present probabilistic incessant spatial [3] [4]-worldly consecutive examples with gap constraints. Such examples are essential for the disclosure of learning given indeterminate direction information. Author propose a dynamic programming approach for processing the recurrence likelihood of these examples, which has direct time intricacy, and Author investigate its inserting into example specification calculations utilizing both broadness first pursuit and profundity first hunt procedures. Our broad experimental study demonstrates the proficiency and viability of our techniques for engineered and real - world datasets.

III. ANALYSIS REPORT

Here analyzed the problem related to pattern matching. Given string T(text) and P(pattern), the pattern matching problem consists of finding a substring of T equal to P. Algorithms are designed and analyzed the problem. [13]. In information retrieval, to improve recall of a web search on a person name or any user entered sting a search engine can automatically expand a query using aliases of the name or string [14]. In our previous example, a user who searches for Amitabh Bacchan might also be interested in retrieving documents in which Bacchan is referred to as Big B. Consequently, we can expand a query on real name using his alias name BigB. The semantic web is intended to solve the entity disambiguation problem by providing a mechanism to add semantic metadata for entities. However, an issue that the semantic web currently faces is that insufficient semantically annotated web contents are available. Automatic extraction of metadata can accelerate the process of semantic annotation. For named entities, automatically extracted aliases can serve as a useful source of metadata, thereby providing a means to disambiguate an entity Revealing the topics inside short messages, for example, tweets and texts, has turned into an essential errand for some content examination applications. Be that as it may, straightforwardly applying customary topic models (e.g. LDA and PLSA) on such short messages may not function admirably. The essential reason lies in that routine topic models verifiably catch the document level word co event examples to uncover topics, and in this manner experience the ill effects of the extreme information sparsity in short records. In this paper, we propose a novel path for demonstrating topics in short messages, alluded as biterm topic model (BTM). In particular, in BTM we take in the topics by specifically displaying the era of word co-event designs (i.e. biterms) in the entire corpus. The real focal topics of BTM are that 1) BTM unequivocally models the word co-event examples to improve the theme learning; and 2) BTM utilizes the accumulated examples as a part of the entire corpus for learning topics to take care of the issue of inadequate word co-event designs at document level. We do broad examinations on real-world short content accumulations. The outcomes exhibit that our approach can find more unmistakable and lucid topics, and fundamentally outflank standard techniques on a few



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 5, Issue 2, February 2017

assessment measurements. Moreover, we find that BTM can beat LDA even on ordinary writings, demonstrating the potential consensus and more extensive utilization of the new point show.

IV. PROPOSED TECHNIQUE

Based on the investigations from the literature survey of various scholars the solution to the problem is achieved which are classified and presented below. In order to characterize user behaviors in published document streams, we study on the correlations among topics extracted from these documents, especially the sequential relations, and specify them as Sequential Topic Patterns (STPs). To solve the innovative and significant problem of mining URSTPs in document streams, many new technical challenges are raised and will be tackled in this paper. Firstly, the input of the task is a textual stream, so existing techniques of sequential pattern mining for probabilistic databases cannot be directly applied to solve this problem. A preprocessing phase is necessary and crucial to get abstract and probabilistic descriptions of documents by topic extraction, and then to recognize complete and repeated activities of Internet users by session identification. Secondly, in view of the real-time requirements in many applications, both the accuracy and the efficiency of mining algorithms are important and should be taken into account, especially for the probability computation process. Thirdly, different from frequent patterns, the user-aware rare pattern concerned here is a new concept and a formal criterion must be well defined, so that it can effectively characterize most of personalized and abnormal behaviors of Internet users, and can adapt to different application scenarios. And correspondingly, unsupervised mining algorithms for this kind of rare patterns need to be designed in a manner different from existing frequent pattern mining algorithms.

1. To the best of our knowledge, this is the first work that gives formal definitions of STPs as well as their rarity measures, and puts forward the problem of mining URSTPs in document streams, in order to characterize and detect personalized and abnormal behaviors of Internet users.
2. A framework to pragmatically solve this problem, and design corresponding algorithms to support it.
3. Preprocessing procedures with heuristic methods for topic extraction and session identification. Then, borrowing the ideas of pattern-growth in uncertain environment, two alternative algorithms are designed to discover all the STP candidates with support values for each user. That provides a trade-off between accuracy and efficiency. At last, we present a user-aware rarity analysis algorithm according to the formally defined criterion to pick out URSTPs and associated users.

Session Identification Since each session should contain a complete publishing behavior of an individual user, we need to at first divide the document stream according to different users, which is an easy job as the author of each document is explicitly given in the input stream. The result for each user u is a subsequence of the topic-level document stream restricted to that user, i.e., $TDS_u = h(td_1, u, t_1), (td_2, u, t_2), \dots, (td_N, u, t_N)$. After that, we also need to partition the subsequence to identify complete and repeated activities as consecutive and non-overlapped sessions. They constitute a session set $S_u = \{s_1, s_2, \dots, s_m\}$ satisfying $TDS_u = s_1 \circ s_2 \circ \dots \circ s_m$, where \circ is the concatenation operator. Time Interval Heuristics. It assumes that the time interval of any two adjacent documents in the same session is less than or equal to a predefined threshold ht_i . The algorithm named $tiPartition$ is shown in Algorithm 1. It examines each document on the input stream orderly to see whether it should be the starting point of a new session, by checking the condition that the time difference between it and its previous document exceeds the threshold Time Span Heuristics. It assumes that the duration of each session is less than or equal to a predefined threshold ht_s . The algorithm is named $tsPartition$, and the only distinction from Algorithm 2 is the condition for a new session in line 3. Specifically, the time point of the previous document t_{i-1} is replaced by the starting time point of the current session t_k , and ht_s is used as threshold. Due to the page limit, the pseudocode is omitted here.

V. CONCLUSION

Mining URSTPs in published document streams on the Internet is a significant and challenging problem. It formulates a new kind of complex event patterns based on document topics, and has wide potential application scenarios, such as real-time monitoring on abnormal behaviors of Internet users. In this paper, several new concepts and the mining problem are formally defined, and a group of algorithms are designed and combined to systematically solve this problem. The experiments conducted on both real (Twitter) and synthetic datasets demonstrate that the proposed approach is very effective and efficient in discovering special users as well as interesting and interpretable URSTPs from Internet document streams, which can well capture users' personalized and abnormal behaviors and



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirccce.com

Vol. 5, Issue 2, February 2017

characteristics. As this paper puts forward an innovative research direction on Web data mining, much work can be built on it in the future

REFERENCES

1. C. C. Aggarwal, Y. Li, J. Wang, and J. Wang, "Frequent pattern mining with uncertain data," in Proc. ACM SIGKDD'09, 2009, pp. 29–38.
2. R. Agrawal and R. Srikant, "Mining sequential patterns," in Proc. IEEE ICDE'95, 1995, pp. 3–14.
3. J. Allan, R. Papka, and V. Lavrenko, "On-line new event detection and tracking," in Proc. ACM SIGIR'98, 1998, pp. 37–45.
4. T. Bernecker, H.-P. Kriegel, M. Renz, F. Verhein, and A. Zuefle, "Probabilistic frequent itemset mining in uncertain databases," in Proc. ACM SIGKDD'09, 2009, pp. 119–128.
5. D. Blei and J. Lafferty, "Correlated topic models," Adv. Neural Inf. Process. Syst., vol. 18, pp. 147–154, 2006.
6. D. M. Blei and J. D. Lafferty, "Dynamic topic models," in Proc. ACM ICML'06, 2006, pp. 113–120.
7. D. Blei, A. Ng, and M. Jordan, "Latent Dirichlet allocation," J. Mach. Learn. Res., vol. 3, pp. 993–1022, 2003.
8. J. Chae, D. Thom, H. Bosch, Y. Jang, R. Maciejewski, D. S. Ebert, and T. Ertl, "Spatiotemporal social media analytics for abnormal event detection and examination using seasonal-trend decomposition," in Proc. IEEE V AST'12, 2012, pp. 143–152.
9. K. Chen, L. Luesukprasert, and S. T. Chou, "Hot topic extraction based on timeline analysis and multidimensional sentence modeling," IEEE Trans. Knowl. Data Eng., vol. 19, no. 8, pp. 1016–1025, 2007.
10. C. K. Chui and B. Kao, "A decremental approach for mining frequent itemsets from uncertain data," in Proc. PAKDD'08, 2008, pp. 64–75.
11. W. Dou, X. Wang, D. Skau, W. Ribarsky, and M. X. Zhou, "LeadLine: Interactive visual analysis of text data through event identification and exploration," in Proc. IEEE V AST'12, 2012, pp. 93–102.
12. G. P. C. Fung, J. X. Yu, P. S. Yu, and H. Lu, "Parameter free bursty events detection in text streams," in Proc. VLDB'05, 2005, pp. 181–192.
13. J. Han, J. Pei, B. Mortazavi-Asl, Q. Chen, U. Dayal, and M. Hsu, "FreeSpan: frequent pattern projected sequential pattern mining," in Proc. ACM SIGKDD'00, 2000, pp. 355–359.
14. N. Hariri, B. Mobasher, and R. Burke, "Context-aware music recommendation based on latent topic sequential patterns," in Proc. ACM RecSys'12, 2012, pp. 131–138.
15. T. Hofmann, "Probabilistic latent semantic indexing," in Proc. ACM SIGIR'99, 1999, pp. 50–57.
16. L. Hong and B. D. Davison, "Empirical study of topic modeling in Twitter," in Proc. ACM SOMA'10, 2010, pp. 80–88.
17. Z. Hu, H. Wang, J. Zhu, M. Li, Y. Qiao, and C. Deng, "Discovery of rare sequential topic patterns in document stream," in Proc. SIAM SDM'14, 2014, pp. 533–541.
18. A. Krause, J. Leskovec, and C. Guestrin, "Data association for topic intensity tracking," in Proc. ACM ICML'06, 2006, pp. 497–504.
19. W. Li and A. McCallum, "Pachinko allocation: DAG-structured mixture models of topic correlations," in Proc. ACM ICML'06, vol. 148, 2006, pp. 577–584.