



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirce.com

Vol. 4, Issue 12, December 2016

A Survey on User-Aware STPs in Document Streams

Somesh D. Kalaskar, Dr. Archana Lomte

M. E. Student, Dept. of Computer, JSPM's BSIOTR, Wagholi, Pune, Maharashtra, India

Asst. Professor, Dept. of Computer, JSPM's BSIOTR, Wagholi, Pune, Maharashtra, India

ABSTRACT: Documents created and distributed on the Internet are ever changing in various forms. Most of existing works are devoted to topic modeling and the evolution of individual topics, while sequential relations of topics in successive documents published by a specific user are ignored. In this paper, in order to characterize and detect personalized and abnormal behaviors of Internet users. Document streams are created and distributed in various forms on the Internet, such as news streams, emails, micro-blog articles, chatting messages, research paper archives, web forum discussions, and so forth. The contents of these documents generally concentrate on some specific topics, which reflect offline social events and users' characteristics in real life. To mine these pieces of information, a lot of researches of text mining focused on extracting topics from document collections and document streams through various probabilistic topic models.

KEYWORDS: UARSTP Recommendation System, Data Mining, Information Retrieval.

I. INTRODUCTION

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



ISSN(Online): 2320-9801
ISSN (Print): 2320-9798

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 4, Issue 12, December 2016

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 (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 PrefixSpan 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 databases [3] From This Paper I Referred-

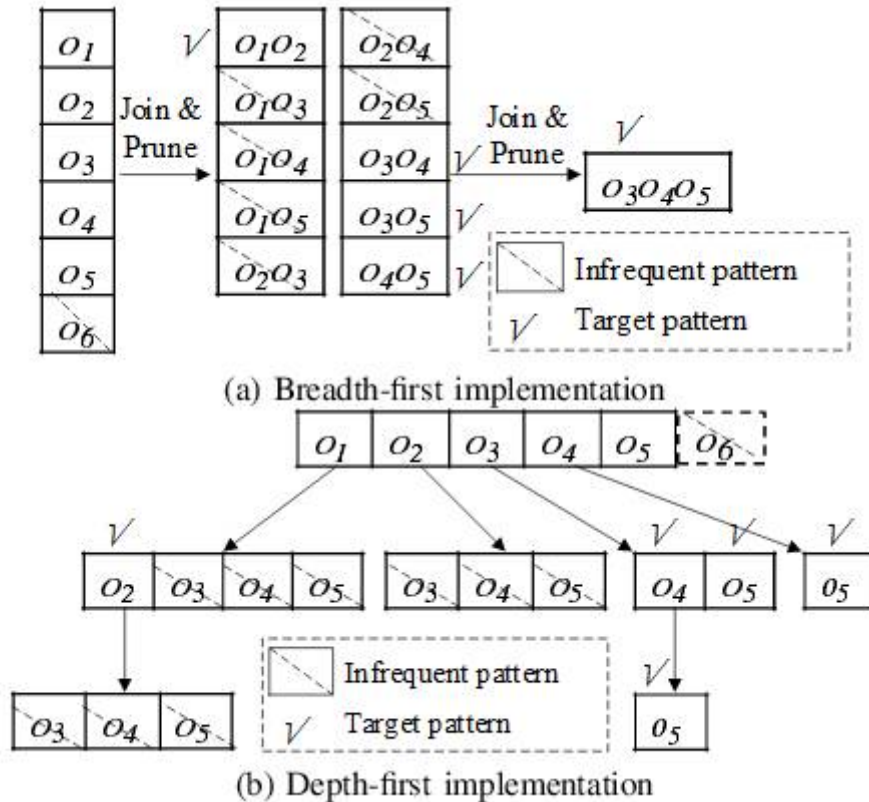
Uncertainty is regular in real - world applications, for instance, in sensor organizes and moving article following , bringing about much enthusiasm for thing set digging for questionable exchange databases. In this paper, author concentrate on example digging for dubious groupings and present probabilistic 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.

International Journal of Innovative Research in Computer and Communication Engineering

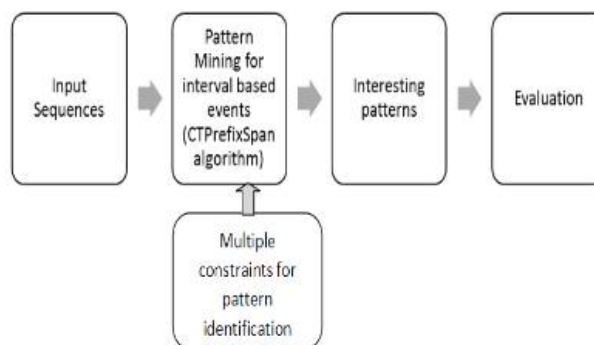
(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 4, Issue 12, December 2016



4. Sequential pattern mining - approaches and algorithms [4] From This Paper IReferred-



Sequence of events, things, or tokens happening in a requested metric space show up regularly in information and the necessity to identify and dissect visit subsequences is a typical issue. Consecutive Pattern Mining emerged as a subfield of information mining to concentrate on this field. This article overviews the methodologies and calculations proposed to date.

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 4, Issue 12, December 2016

5. A biterm topic model for short texts[5] From This Paper IReferred-

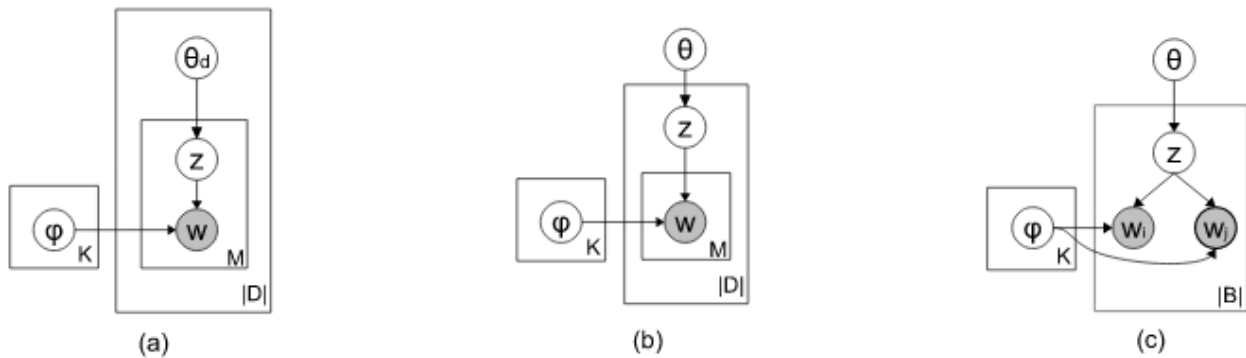


Fig.3: Graphical representation of (a) LDA, (b) mixture of unigrams, and (c) BTM

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 sparsely in short records. In this paper [5], Authorproposes 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. Author 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 assessment measurements. Moreover, author find that BTM can beat LDA even on ordinary writings, demonstrating the potential consensus and more extensive utilization of the new point show.

III. SYSTEM ARCHITECTURE

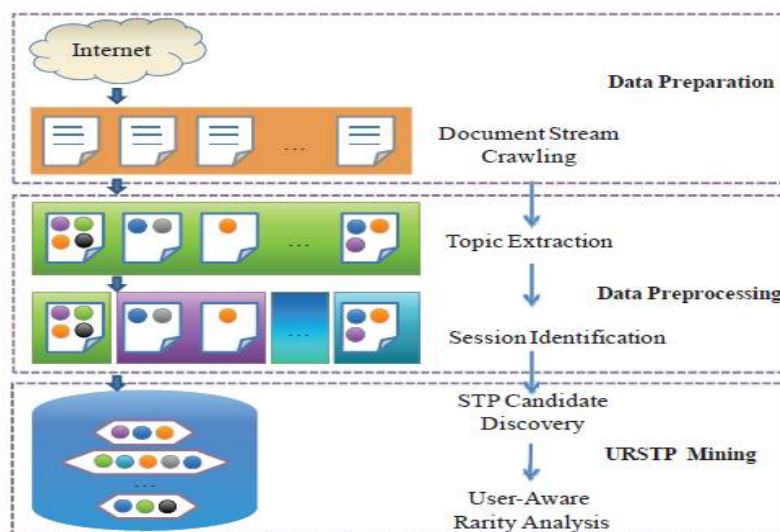


Fig4. System architecture



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijirce.com

Vol. 4, Issue 12, December 2016

In order to characterize and detect personalized and abnormal behaviors of Internet users, we propose Sequential Topic Patterns (STPs) and formulate the problem of mining User-aware Rare Sequential Topic Patterns (URSTPs) in document streams on the Internet. 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). Each of them records the complete and repeated behavior of a user when she is publishing a series. Topic mining in document collections has been extensively studied in the literature. Topic Detection and Tracking (TDT) task aimed to detect and track topics (events) in news streams with clustering-based techniques on keywords. 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 characteristics.

Taking advantage of these extracted topics in document streams, most of exist works analysed the evolution of individual topics to detect and predict social events as well as user behaviours. In order to find significant STPs, a document stream should be divided into independent sessions in advance with the definition. A sketch map of session identification each ellipse represents a session, and all the sessions in each line constitute a document subsequence for a specific user. A can conclude that the two algorithms have their respective advantages. Which one is appropriate for the real task reflects a trade-off between mining accuracy and execution speed, and should depend on the specific requirements of application scenarios.

IV. 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 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] 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. International Journal of Multimedia Information Retrieval, 2014, 3.1: 29 - 39.
- [2] Z. Zhao, D. Yan, and W. Ng, "Mining probabilistically frequent sequential patterns in large uncertain databases," IEEE Trans. Knowl. Data Eng., vol. 26, no. 5, pp. 1171 – 1184, 2014.
- [3] Y. Li, J. Bailey, L. Kulik, and J. Pei, "Mining probabilistic frequent spatio-temporal sequential patterns with gap constraints from uncertain databases," in Proc. IEEE ICDM'13, 2013, pp. 448 – 457.
- [4] C. H. Mooney and J. F. Roddick, "Sequential pattern mining - approaches and algorithms," ACM Comput. Surv., vol. 45, no. 2, pp. 19:1 – 19:39, 2013.
- [5] 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,"
- [6] In Proc. IEEE VAST'12, 2012, pp. 143–152.
- [7] X. Yan, J. Guo, Y. L an, and X. Cheng, "A bitern topic model for short texts," in Proc. ACM WWW'13, 2013, pp. 1445 – 1456.
- [8] 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.
- [9] C. K. Chui and B. Kao, "A decremental approach for mining frequent itemsets from uncertain data," in Proc. PAKDD'08, 2008, pp. 64–75.
- [10] 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 VAST'12, 2012, pp. 93–102.
- [11] 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.



ISSN(Online): 2320-9801
ISSN (Print): 2320-9798

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: www.ijircce.com

Vol. 4, Issue 12, December 2016

BIOGRAPHY

Mr. Somesh D.Kalaskar: is pursuing M.E. at Computer Dept. JSPM's BhivarabaiSawant Institute of Technology and Research and Technology, PuneMaharashtra, India and have completed B.Tech from Walchand College of Engineering Sangli. His area of interest is data mining.

Dr. ArchanaLomte is an Assistant Professor at, Department of Computer Engineering, JSPM's BhivarabaiSawant Institute of Technology and Research, Pune, Maharashtra, India.