

Frequent Itemset Mining: A Metadata Based Approach for Knowledge Discovery

Basavaraj A.^{1*} Goudannavar², Prashant Bhat³

^{1*,2}Department of Computer Science, Rani Channamma University, Belagavi-591156, Karnataka, India

²Department of Business Analytics/Data Science, Chris Institute of Management, Lavasa-412112, Pune, Maharashtra, India

Available online at: www.ijcseonline.org

Received: 17/Feb//2018, Revised: 26/Feb2018, Accepted: 19/Mar/2018, Published: 30/Mar/2018

Abstract— Frequent sets play a crucial role in many Data Mining tasks that try to find interesting patterns from databases, such as correlations, association rules, classification and clustering. The Association Rules is one of the most used functions in data mining. The method is used both database researchers and data mining users. In this article, association rule mining algorithms are discussed and demonstrated. Mining Associate rule algorithm that search for approximate strong association rules from multimedia databases. The Apriori-like sequential pattern mining approach based on candidate generates-and test can also be explored by mapping a sequence multimedia database into vertical data format. This approach is useful to finding frequent itemsets, which probabilistic frequent itemsets based on possible datasets.

Keywords— Web Multimedia Mining, Association rule, Frequent itemsets, Knowledge discovery

I. INTRODUCTION

Association rule mining is a data mining task that discovers relationships among items in a transactional database. Association rules have been extensively studied in the literature for their usefulness in many application domains such as diagnosis decisions support, telecommunication, multimedia database, etc. The efficient discovery of such rules has been a major focus in the data mining research community. From the original *apriori* algorithm [1] there have been a remarkable number of variants and improvements of association rule mining algorithms [2].

Association rules are required to satisfy both a minimum support and a minimum confidence constraint at the same time. At medium to low support values, often a great number of frequent itemsets are found in a database. However, since the definition of support enforces that all subsets of a frequent itemset have to be also frequent, it is sufficient to only mine all maximal frequent itemsets, defined as frequent itemsets which are not proper subsets of any other frequent itemset. Another approach to reduce the number of mined itemsets is to only mine frequent closed itemsets. An itemset is closed if no proper superset of the itemset is contained in each transaction in which the itemset is contained [3]. In this proposed work, using KNIME data mining tool is used [4], the web multimedia-video metadata are extracted and frequent itemsets and closed frequent itemsets on available metadata of web multimedia-videos using Apriori algorithm and FP-Growth algorithms. The frequent pattern mining results are analyzed.

The respite of the paper is ordered as follows: The part 2 represents related works on the association rule on web multimedia datasets, part 3 represents proposed

methodology, part 4 represents evaluation analysis of association rule, and finally part 5 represents conclusion.

II. PRIOR WORKS

The authors [5] introduce an improvement to the mining Apriori association rule generates approximate association rules. The apriori algorithm takes into consideration missing values and noisy data. The experiment indicates that apriori effectively generates rules that approximate positive correlations in the input database.

In the proposed work, multiple database passes to generate the frequent item sets [6]. In this article a new Algorithm is proposed which can mine frequent patterns with a single scan of database. The time taken by this algorithm is compared with the other algorithms also and it proves that time taken by this algorithm is less than the partition algorithm and the Apriori but more than the time taken by FP-Tree Algorithms but FP tree has its own disadvantages such as if every transaction in database is different then we will have as many as $2n$ leaves in FP tree (where n is the no. of items) [7] [8]. In this article, proposed four parallel versions of a novel sequential mining algorithm for discovery of frequent itemsets are proposed. The parallel algorithms were compared analytically and experimentally, with respect to some factors, such as communication rate, response time, computation/communication ratio and load balancing. [9].

A novel adaptive classification method using random forests, which is a machine learning algorithm with proven good performance on many traditional classification problems [10]. Video classification is the first step toward

multimedia content understanding. When video is classified into conceptual categories, it is usually desirable to combine evidence from multiple modalities [11].

III. PROPOSED METHODOLOGY

In this part we describe an effectual methodology to extract the metadata from web multimedia files and categorize them based on the frequent itemsets of metadata

by applying data mining techniques. In this proposed method, out of the total metadata dataset, 60% are used for training and remaining 40% are used for testing the classification model built using Decision Tree and SVM classification methods. The consequences are analyzed and the effectiveness of the proposed method has been demonstrated.

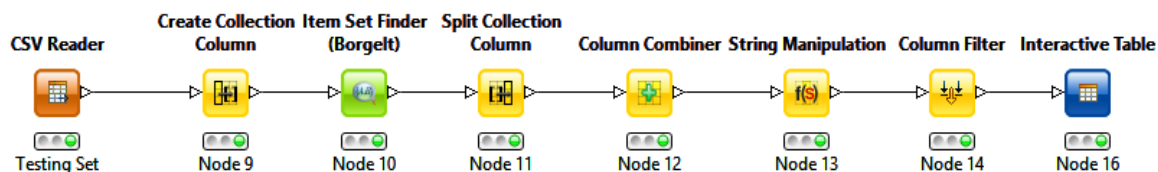


Figure 1: System model of the proposed methodology

The proposed system model is represented in Figure 1. It consists of the following components:

- i) Web multimedia-video metadata extraction and Pre-processing
- ii) Frequent pattern mining and association rules model
- iii) Result analysis

The functionality of each element of the proposed system model is discussed in the following subsections.

3.1 Web Multimedia-Video Metadata Extraction and pre-processing

There are many tools available to decompose the web multimedia data meaning that, it is possible to retrieve individual components object from multimedia data. The various techniques to extract web multimedia data and pre-processing are discussed in our previous work [12].

3.2 Frequent pattern mining and association rules model

The proposed work, we adopt a model to mine frequent itemsets of web multimedia metadata. The mining precision and efficiency will depend on the constructed frequent itemset mining model. This section represents detailed procedure to construct frequent itemset model.

3.2.1 Association rules

Association rules are created by analyzing data for frequent if/then patterns and using the criteria support and confidence to identify the most important relationships. Support is an indication of how frequently the items appear in the database. Confidence indicates the number of times

4. Is the candidate set empty, if not goto 2
5. For each frequent itemset I , generate all nonempty subsets of I

the if/then statements have been found to be true. Given a set of transactions, each described by an unordered set of items, an association rule $A \rightarrow B$ may be discovered in the data, where A and B are conjunctions of items. The intuitive meaning of such a rule is that transactions in the database which contain the items in A , tend to also contain the items in B . An example of such a rule might be that observed metadata multimedia datasets to find the frequent itemset mining also association rule generating. In this case, $A = \{\text{Stream size Mib, Image Height, Audio duration}\}$ and $B = \{\text{Video duration}\}$. Two numbers are associated with each rule that indicate the support and confidence of the rule using multimedia data.

The problem of discovering association rules can be divided into two steps:

1. Find all *itemsets* (sets of items appearing together in a transaction) whose support is greater than the specified threshold.

2. Generate association rules from the frequent itemsets.

Confidence of a candidate rule is calculated as support of the multimedia metadata. All rules that meet the confidence threshold are reported as discoveries of the algorithm.

1. Scan the (entire) transaction database to get the support S of each 1-itemset, compare S with min_sup , and get a set of frequent 1-itemsets, L_1
2. Use L_{k-1} join L_{k-1} to generate a set of candidate kitemsets.
3. Scan the transaction database to get the support S of each candidate k -itemset in the final set, compare S with min_sup , and get a set of frequent kitemsets, L_k
6. For every nonempty subset s of I , output the rule $S = (I-s)$ if its confidence $C > \text{min} > \text{conf}$

Algorithm 1: Apriori Algorithm

In this algorithm first calculates single item frequencies to determine the frequent 1-itemsets (Minimum support of Multimedia data). Once candidates are generated, itemsets are removed from consideration if any (k-1) subset of the candidate is not in Lk-1.

4. Experimental Results and Discussions

4.1 Association rule and frequent itemset analysis

The association rule algorithm used in this experiment is restricted to finding rule and frequent itemset between pairs of items only. In this experiment propose a novel approach for frequent item sets mining using multimedia metadata. Frequent itemsets is an item set that satisfies minimum support. The multimedia metadata information fields like video duration, video bit rate kbps, maximum bit rate kbps, width pixels, height pixels, display aspect ratio, bits/(pixel*frame), stream size mib, audio duration, audio bit rate kbps, maximum bit rate kbps, stream size mib, image resolution, image height, image width, text page, word count, character count, line count, paragraph count, size in kbps, class. The purpose of this experiment is to find out the result for generating frequent itemsets using Apriori algorithm, and added on the item that it took to find the association rules within the frequent itemsets.

The first step of Apriori is to count the frequencies in data set, called the supports, of each metadata item separately.

Table-1: Generation of frequent 1-itemset

Row ID	ItemSetSize	ItemSetSupport	RelativeItemSetSupport%	combined string
Row0	1	493	100	"1", "2", "3", "4", "5", "6", "7", "8", "9", "10", "11", "12", "13", "14", "15", "16", "17", "18", "19", "20", "21", "22", "23", "24", "25", "26", "27", "28", "29", "30", "31", "32", "33", "34", "35", "36", "37", "38", "39", "40", "41", "42", "43", "44", "45", "46", "47", "48", "49", "50", "51", "52", "53", "54", "55", "56", "57", "58", "59", "60", "61", "62", "63", "64", "65", "66", "67", "68", "69", "70", "71", "72", "73", "74", "75", "76", "77", "78", "79", "80", "81", "82", "83", "84", "85", "86", "87", "88", "89", "90", "91", "92", "93", "94", "95", "96", "97", "98", "99", "100", "101", "102", "103", "104", "105", "106", "107", "108", "109", "110", "111", "112", "113", "114", "115", "116", "117", "118", "119", "120", "121", "122", "123", "124", "125", "126", "127", "128", "129", "130", "131", "132", "133", "134", "135", "136", "137", "138", "139", "140", "141", "142", "143", "144", "145", "146", "147", "148", "149", "150", "151", "152", "153", "154", "155", "156", "157", "158", "159", "160", "161", "162", "163", "164", "165", "166", "167", "168", "169", "170", "171", "172", "173", "174", "175", "176", "177", "178", "179", "180", "181", "182", "183", "184", "185", "186", "187", "188", "189", "190", "191", "192", "193", "194", "195", "196", "197", "198", "199", "200", "201", "202", "203", "204", "205", "206", "207", "208", "209", "210", "211", "212", "213", "214", "215", "216", "217", "218", "219", "220", "221", "222", "223", "224", "225", "226", "227", "228", "229", "230", "231", "232", "233", "234", "235", "236", "237", "238", "239", "240", "241", "242", "243", "244", "245", "246", "247", "248", "249", "250", "251", "252", "253", "254", "255", "256", "257", "258", "259", "260", "261", "262", "263", "264", "265", "266", "267", "268", "269", "270", "271", "272", "273", "274", "275", "276", "277", "278", "279", "280", "281", "282", "283", "284", "285", "286", "287", "288", "289", "290", "291", "292", "293", "294", "295", "296", "297", "298", "299", "300", "301", "302", "303", "304", "305", "306", "307", "308", "309", "310", "311", "312", "313", "314", "315", "316", "317", "318", "319", "320", "321", "322", "323", "324", "325", "326", "327", "328", "329", "330", "331", "332", "333", "334", "335", "336", "337", "338", "339", "340", "341", "342", "343", "344", "345", "346", "347", "348", "349", "350", "351", "352", "353", "354", "355", "356", "357", "358", "359", "360", "361", "362", "363", "364", "365", "366", "367", "368", "369", "370", "371", "372", "373", "374", "375", "376", "377", "378", "379", "380", "381", "382", "383", "384", "385", "386", "387", "388", "389", "390", "391", "392", "393", "394", "395", "396", "397", "398", "399", "400", "401", "402", "403", "404", "405", "406", "407", "408", "409", "410", "411", "412", "413", "414", "415", "416", "417", "418", "419", "420", "421", "422", "423", "424", "425", "426", "427", "428", "429", "430", "431", "432", "433", "434", "435", "436", "437", "438", "439", "440", "441", "442", "443", "444", "445", "446", "447", "448", "449", "450", "451", "452", "453", "454", "455", "456", "457", "458", "459", "460", "461", "462", "463", "464", "465", "466", "467", "468", "469", "470", "471", "472", "473", "474", "475", "476", "477", "478", "479", "480", "481", "482", "483", "484", "485", "486", "487", "488", "489", "490", "491", "492", "493", "494", "495", "496", "497", "498", "499", "500"

The support $supp(A)$ of an itemset A is defined as the proportion of transactions in the dataset which contain the itemset i.e.,

$$supp(A) = \frac{\text{no. of metadata which contain the itemset } A}{\text{total no. of metadata}}$$

In the multimedia metadata, each item is a member of the set of candidate 1-itemsets, C_1 . The algorithm simply scans all of the transactions in order to count the number of occurrences of each item. The minimum support count required is 99, that is, $min\ sup = 99$. (Here, we are referring to absolute support because we are using a support count. The corresponding relative support is $99/493 = 20.08\%$). To be even more explicit the experiment multimedia metadata is 99 is the number of transactions from the metadata which contain the itemset while 493 represents the total number of transactions. We can define a minimum support level to qualify as "frequent," which depends on the context. For this case, let minimum support 99 and relative support are 20.08. Therefore, all are frequent in 1-itemset.

The next step is to generate a list of all 2-pairs of the frequent items. Had any of the above items not been frequent, they would not have been included as a possible member of possible 2-item pairs. In this way, Apriori prunes the tree of all possible data sets. Further we again select only these items (now 2-pairs are items). The web multimedia metadata 205 transaction database in 1-itemset, we found 186 which are frequent in 2-itemset, the remaining 19 items are rejected which are lower than the support value are shown in table-2.

Table 2: Generation of frequent 2-itemset

Row ID	ItemSetSize	ItemSetSupport	RelativeItemSetSupport%	combined string
Row0	2	102	20.69	"1280", "1", "2", "3", "4", "5", "6", "7", "8", "9", "10", "11", "12", "13", "14", "15", "16", "17", "18", "19", "20", "21", "22", "23", "24", "25", "26", "27", "28", "29", "30", "31", "32", "33", "34", "35", "36", "37", "38", "39", "40", "41", "42", "43", "44", "45", "46", "47", "48", "49", "50", "51", "52", "53", "54", "55", "56", "57", "58", "59", "60", "61", "62", "63", "64", "65", "66", "67", "68", "69", "70", "71", "72", "73", "74", "75", "76", "77", "78", "79", "80", "81", "82", "83", "84", "85", "86", "87", "88", "89", "90", "91", "92", "93", "94", "95", "96", "97", "98", "99", "100", "101", "102", "103", "104", "105", "106", "107", "108", "109", "110", "111", "112", "113", "114", "115", "116", "117", "118", "119", "120", "121", "122", "123", "124", "125", "126", "127", "128", "129", "130", "131", "132", "133", "134", "135", "136", "137", "138", "139", "140", "141", "142", "143", "144", "145", "146", "147", "148", "149", "150", "151", "152", "153", "154", "155", "156", "157", "158", "159", "160", "161", "162", "163", "164", "165", "166", "167", "168", "169", "170", "171", "172", "173", "174", "175", "176", "177", "178", "179", "180", "181", "182", "183", "184", "185", "186", "187", "188", "189", "190", "191", "192", "193", "194", "195", "196", "197", "198", "199", "200", "201", "202", "203", "204", "205", "206", "207", "208", "209", "210", "211", "212", "213", "214", "215", "216", "217", "218", "219", "220", "221", "222", "223", "224", "225", "226", "227", "228", "229", "230", "231", "232", "233", "234", "235", "236", "237", "238", "239", "240", "241", "242", "243", "244", "245", "246", "247", "248", "249", "250", "251", "252", "253", "254", "255", "256", "257", "258", "259", "260", "261", "262", "263", "264", "265", "266", "267", "268", "269", "270", "271", "272", "273", "274", "275", "276", "277", "278", "279", "280", "281", "282", "283", "284", "285", "286", "287", "288", "289", "290", "291", "292", "293", "294", "295", "296", "297", "298", "299", "300", "301", "302", "303", "304", "305", "306", "307", "308", "309", "310", "311", "312", "313", "314", "315", "316", "317", "318", "319", "320", "321", "322", "323", "324", "325", "326", "327", "328", "329", "330", "331", "332", "333", "334", "335", "336", "337", "338", "339", "340", "341", "342", "343", "344", "345", "346", "347", "348", "349", "350", "351", "352", "353", "354", "355", "356", "357", "358", "359", "360", "361", "362", "363", "364", "365", "366", "367", "368", "369", "370", "371", "372", "373", "374", "375", "376", "377", "378", "379", "380", "381", "382", "383", "384", "385", "386", "387", "388", "389", "390", "391", "392", "393", "394", "395", "396", "397", "398", "399", "400", "401", "402", "403", "404", "405", "406", "407", "408", "409", "410", "411", "412", "413", "414", "415", "416", "417", "418", "419", "420", "421", "422", "423", "424", "425", "426", "427", "428", "429", "430", "431", "432", "433", "434", "435", "436", "437", "438", "439", "440", "441", "442", "443", "444", "445", "446", "447", "448", "449", "450", "451", "452", "453", "454", "455", "456", "457", "458", "459", "460", "461", "462", "463", "464", "465", "466", "467", "468", "469", "470", "471", "472", "473", "474", "475", "476", "477", "478", "479", "480", "481", "482", "483", "484", "485", "486", "487", "488", "489", "490", "491", "492", "493", "494", "495", "496", "497", "498", "499", "500"

Generation of all 3-triples of list the frequent items from the dataset. The above items not been frequent, they would not have been included as a possible member of possible 3-item pairs. Further we again select only these items (now 3-pairs are items) the 186 transaction database in 2-itemset, we found 136 transaction dataset which are frequent in 3-itemset, the remaining 50 items are rejected which are lower than the support value are shown in table-3.

Confidence = $99/154 = 64\%$

Confidence = $99/143 = 69\%$

If the minimum confidence threshold is, say, 20%, All of these association rules are considered strong because they meet the minimum confidence and support which was set at 20%.

5. Conclusion

In this paper, Apriori is one of the most well data mining approaches is to find frequent itemsets from a transaction dataset of multimedia metadata and to generate strong association rules. Once frequent itemsets are obtained, it is straightforward to generate association rules with confidence larger than or equal to a user specified minimum confidence. Further we use apriori algorithms for finding informative patterns from the complex multimedia data sources.

REFERENCES

- [1] Agrawal, R., Imielinski, T., Swami, "A.: Mining association rules between sets of items in large databases",
- [2] Acm SIGKDD Explorations Newsletter, Volume 22, Issue 2, June 1993, ISBN: 0-89791-592-5.
- [3] Bart Goethals, Mohammed J. Zaki, "Advances in Frequent Itemset Mining Implementations Report on FIMI'03", ACM SIGKDD Explorations Newsletter, Volume 6, Issue 1, June 2004, ISSN: 1931-0145.
- [4] Michael Hahsler, Bettina Grun and Kurt Hornik, "arules – A Computational Environment for Mining Association Rules and Frequent Item Sets", Journal of Statistical Software, Volume 14, Issue 15, October 2005,.
- [5] ChenGang "MediaInfo extractor – A Tool for Media Data Mining", 2011. <http://mediaarea.net/en/MediaInfo>.
- [6] Jyothsna R. Nayak and Diane J. Cook, "Approximate Association Rule Mining", FLAIRS-01 Proceedings. 2001, AAAI (www.aaai.org).
- [7] Syed Khairuzzaman Tanbeer , Chowdhury , Farhan Ahmed and Byeong-Soo Jeong , "Parallel and Distributed Algorithms for FP mining in large Databases", IETE Technical Review Vol 26, Issue 1 , pp 55-65, Jan 2009.
- [8] Pradeep Chouksey* Juhi Singh, R.S. Thakur and R.C. Jain, "Frequent Pattern Mining using Candidate Generation approach with Single Scan of Database", Symposium on Progress in Information & Communication Technology 2009
- [9] Ms. Manali Rajeev Raut, Ms. Hemlata Dakhore, "Association Rule Mining in Horizontally Distributed Databases", Manali Rajeev Raut et al, / (IJCSIT) International Journal of Computer Science and Information Technologies, Vol. 5 (6) , 2014, 7540-7544
- [10] Zhi Liu, Tianhong Sun and Guoming Sang, " An Algorithm of Association Rules Mining in Large Databases Based on Sampling ", International Journal of Database Theory and Application Vol.6, No.6 , 2013.
- [11] Renáta Iváncsy, István Vajk', "Frequent Pattern Mining in Web Log Data", Acta Polytechnica Hungarica Vol. 3, No. 1, 2006.
- [12] S. M. Fakhrahmad1 And Gh. Dastghaibyfar2, "An Efficient Frequent Pattern Mining Method and its Parallelization in Transactional Databases", Journal Of Information Science And Engineering 27, 511-525 (2011).
- [13] Siddu P. Algur1, Basavaraj A. Goudannavar2*, "Web Multimedia Mining: Metadata Based Classification and Analysis of Web Multimedia", International Journal of Advanced Research in Computer Science and Software Engineering (IJARCSSE) November-2015, pp. 324-330 ISSN: 2277-128X

Authors

Dr. Basavaraj A. Goudannavar is working as Lecturer, Dept. of Computer Science, Rani Channamma University, Belagavi, Karnataka, India. He received BCA and MCA degrees from Karnatak University, Dharwad, Karnataka, India, in 2005 and 2008 respectively. His research interest includes Data Mining, Web Mining, Web multimedia mining, and Knowledge discovery techniques. He published 17 research papers in International Journals and International conferences. He has attended and participated in International and National Conferences and Workshops in his research field.



Dr. Prashant Bhat is working as Assistant Professor, Department of Business Analytics/Data Science, Chris Institute of Management, Lavasa- 412112, Pune, Maharashtra, India. He received B.Sc and M.Sc (Computer Science) degrees from Karnatak University, Dharwad, Karnataka, India, in 2010 and 2012 respectively. His research interest includes Data Mining, Web Mining, web multimedia mining and Information Retrieval from the web and Knowledge discovery techniques, and published 22 research papers in International Journals. Also he has attended and participated in International and National Conferences and Workshops in his research field.

