

TDMCS: An Efficient Method for Mining Closed Frequent Patterns over Data Streams Based on Time Decay Model

Meng Han, Jian Ding, and Juan Li

School of Computer Science and Engineering, Beifang University of Nationalities, China

Abstract: In some data stream applications, the information embedded in the data arriving in the new recent time period is important than historical transactions. Because data stream is changing over time, concept drift problem will appear in data stream mining. Frequent pattern mining always generate useless and redundant patterns, in order to obtain the result set of lossless compression, closed pattern is needed. A novel method for efficiently mining closed frequent patterns on data stream is proposed in this paper. The main work includes: distinguish importance of recent transactions from historical transactions based on time decay model and sliding window model; design and use frame minimum support count-maximal support error threshold-decay factor (θ - ϵ - f) to avoid concept drift; use closure operator to improve the efficiency of algorithm; design a novel way to set decay factor: average-decay-factor faverage in order to balance the high recall and high precision of algorithm. The performance of proposed method is evaluated via experiments, and the results show that the proposed method is efficient and steady-state, it applies to mine data streams with high density and long patterns, it is suitable for different size sliding windows, and it is also superior to other analogous algorithms.

Keywords: data stream mining, frequent pattern mining, closed pattern mining, time decay model, sliding window, concept drift.

Received January 15, 2015; accepted August 12, 2015

1. Introduction¹

Data stream as a new data model is widely used in many applications. Data stream which is different from traditional database is time ordered, rapidly changing, massive and unlimited. Searching for frequent patterns in a continuous data stream has become important and challenging.

In recent years, some algorithms for mining frequent patterns or itemsets on data streams have been proposed. Algorithms such as Sticky Sampling[15], Lossy Counting[15], XSM[1] and EDPM[26] mine frequent patterns which meet maximal support error rate and minimum support count. These methods do not distinguish between recent and historical transactions and do not consider the importance of recent transactions. In addition, these methods for mining complete result sets will produce a lot of useless patterns. For reducing the number of patterns, concise pattern set should be mined, mainly including: maximal frequent patterns, closed frequent patterns, top- k frequent patterns or a combination of them and so on. Algorithms Max-FISM[6] and GUIDE[19] discover recent maximal frequent patterns based on sliding windows. WMFP-SW[10] mines weighted maximal frequent patterns based on sliding windows. Algorithms Moment[5], NewMoment[11], CloStream[24],

Stream_FCI[20], TMoment[17], IncMine[4] and CloStream*[25] discover closed frequent patterns based on sliding windows. TOPSIL- Miner[23] uses landmark windows to mine top- k frequent patterns. Methods Top- k Lossy Counting[22], MSWTP[2] and Top- k Miner[18] discover top- k frequent patterns based on sliding windows. FCI_max[21] mines closed top- k frequent patterns based on sliding windows and so on. The drawbacks of above algorithms are that: (1) using only the minimum support threshold for frequent patterns mining and unprocessed concept drift problem of data streams. (2) Although window model are used in these methods, the weights of transactions in window are set to the same weights.

As can be seen from the above algorithms, mining frequent patterns on data streams usually based on window model, especially the sliding window. The reason is that recent transactions normally contain more information than historical ones. Besides sliding window model, Time Decay Model(TDM)[3, 5, 8-9, 12-14,19] is also used to process recent transactions.

TDM-based methods to mine frequent patterns on data stream emphasize that the importances of recent and historical transactions should be distinguished in the window. Recent years, the ways to set decay factor in TDM usually divide into two categories. The first one set decay factor to random value in the range of (0, 1) [9, 14, 19]. Such ways lead to the instability of the mining results because of the random values of decay factor. The second method assumes that

¹ This paper is supported by National Nature Science Foundation of China (61563001), Science Foundation of State Nationalities Affairs Commission (14BFZ008) and Beifang University of Nationalities (2014XYZ13).

algorithm meets 100% Recall or 100% Precision to get the upper and lower bounds of decay factors[3, 12]. Then set decay factor to the upper bound or lower bound or random value between them. The problems of these two ways to set decay factor are that it can get high recall or high precision of algorithm, while get the low corresponding precision or recall of algorithm. Or because of the uncertainty of the decay factor value, the result sets of algorithm are instability.

In order to avoid concept drift, distinguish recent transactions from historical ones, discover compact pattern result set efficiently, and apply to mine high dense transactions and long patterns, a novel algorithm is proposed in this paper. Mainly works and innovations are that: (1) design a novel way to set decay factor f . Existed methods set f to boundary value of lower bound and high bound by assuming 100% Recall and 100% Precision[3, 12, 26], or to a random value in range of (0, 1)[13, 16]. The former will lead to corresponding algorithm low Precision or low Recall. And the later will make unstable performance of algorithm. In order to balance Recall and Precision of algorithm, proposed an average way to set decay factor in this paper. (2) Three layers frame: minimum support-maximum support error-decay factor is used in this paper to solve the concept drift problem and avoid loss of possible frequent patterns. (3) Proposed a novel algorithm to mine closed frequent patterns on data streams based on time decay model and sliding window model. It can get lossless compression result set. Time decay model[3, 12-13, 16] is used to further emphasize the importance of recent transactions and reduce the importance of historical ones. By the comparison of precisions of novel algorithm and existed algorithms, the novel algorithm can get more accurate result sets.

The rest of this paper is organized as follows. Section 2 presents background knowledge; mainly on closure operator and time decay model. The efficient novel algorithm based on time decay model to discover closed patterns is detailed in section 3. Section 4 describes the experiments and explains the experimental results. Section 5 concludes this work.

2. Preliminaries

A data stream $DS = \langle T_1, T_2, \dots, T_i, \dots \rangle$ is a continuous and unbounded sequence of transactions in a timely order, where T_i ($i=1, 2, \dots$) is the i th transaction. Each transaction contains a unique transaction identifier t_{id} , as shown in the first column of Table 1. The support count of frequent pattern P , denoted as $freq(P, N)$ [6], is the number of transactions in existed N transactions in which P occurs.

- **Define 1.** (Frequent Pattern[12]) Let N be the sliding window size, and θ ($\theta \in (0,1]$) be the minimum support. If itemset P meets $freq(P, N) \geq \theta \times N$, P is frequent pattern.

- **Define 2.** (Half-Frequent Pattern, Non-Frequent Pattern[12]) Let N be the sliding window size, θ ($\theta \in (0,1]$) be the minimum support and ε ($\varepsilon \in (0, \theta)$) be the maximal support error. If itemset P meets $\theta \times N \geq freq(P, N) \geq \varepsilon \times \theta \times N$, P is half-frequent pattern. Else if $freq(P, N) < \varepsilon \times \theta \times N$, P is non-frequent pattern.

Table 1 Transaction data stream

TID	Transaction
t_1	1 3 4
t_2	2 3 5
t_3	1 2 3 5
t_4	2 3 4 5

Data stream changes in real time and the infrequent patterns over time may become frequent patterns. That is concept drift. Therefore, in order to reduce the number of missing possible patterns, frequent patterns and half-frequent patterns need to be maintained during mining process. In addition, in order to reduce the cost of maintaining patterns, non-frequent patterns need to be lost. By this way, the possible error of loss patterns is not greater than ε [3, 12]. Therefore, using θ - ε framework can solve the problem of concept drift.

A heavy problem of mining frequent pattern from data stream is generated a large number of useless patterns. Therefore, mining useful and compressed pattern is needed. Discovering closed frequent pattern is a common method, which is lossless compressed and contains all the information of the complete result. Meanwhile, in order to improve the efficiency to discover closed patterns, closure operator[24-25] is used in this paper. The performance of the algorithm with closure operator is better than classic closed pattern mining algorithms such as Moment[5], CFI-Stream[9] and NewMoment[11]. Take closure operator into account, the concept of closed patterns are shown in definitions 3 to 5.

- **Define 3.** (Closure Operator [24-25]) Let T be the subsets of all that transactions in D , denotes as $T \subseteq D$. Let Y be the subsets of all items set I ($Y \subseteq I$) which appears in D . Concept of closed itemset is based on the following two functions h and g :

$$h(T) = \{i \in I \mid \forall t \in T, i \in t\} \quad (1)$$

$$g(Y) = \{t \in D \mid \forall i \in Y, i \in t\} \quad (2)$$

Function h takes T as input and returns an itemset included in all transactions belonging to T . Function g takes an itemset Y as an input and returns a set of transactions including Y . A function

$$C = h \circ g = h(g)$$

is called Closure Operator.

- **Define 4.** (Closed Itemset[25]) An itemset P is called a closed itemset if and only if it satisfy

Formula 3. Otherwise, P is non-closed. The $C(P)$ is called the closure of P .

$$C(P) = h(g(P)) = P \quad (3)$$

- **Define 5.** (Closed Frequent/Half-Frequent Pattern) If itemset $P=C(P)$ and its support is no less than minimum support, then P is called a closed frequent pattern. If itemset $P=C(P)$ and its support is no less than maximal support error, then P is called a closed half-frequent pattern. Otherwise, P is non-frequent pattern.

Due to the continuous and infinite, knowledge contained in data stream will change with the passage of time. Under normal circumstances, the value of recent transaction is important than historical one. Therefore, it is necessary to increase the weight of recent transaction. A time decay model (TDM) is developed to gradually decay the occurrence count of itemset contained in the transaction[2, 11]. Let the decay ratio of support count in the unit time is decay factor f ($f \in (0,1]$). When T_n arrives, support count of frequent pattern P is denoted as $freq_d(P, T_n)$. Each time a new transaction arrives, $freq_d(P, T_n)$ is multiplied by a decay factor f . When the m th transaction T_m arrives, r is 1 if it contains P , otherwise $r=0$. The $freq_d(P, T_m)$ based on decay factor is shown in Formulas 4 and 5.

$$\begin{aligned} freq_d(P, T_m) &= r, \quad \text{if } m=1 \\ freq_d(P, T_m) &= freq_d(P, T_{m-1}) * f + r, \quad \text{if } m \geq 2 \end{aligned} \quad (4)$$

$$\begin{aligned} r &= 1, \quad \text{if } P \subseteq T_m \\ r &= 0, \quad \text{otherwise} \end{aligned} \quad (5)$$

3. Algorithm TDMCS

In this section, data structures and the new way to define decay factor f are introduced, and the proposed algorithm TDMCS (TDM-Based Closed Frequent Pattern Mining on Data Stream) is introduced in detail which is used to mining frequent closed patterns based on f - θ - ϵ framework.

Three data structures are used in algorithm TDMCS, including *ClosedTable*[24], *CidList*[24] and *NewTransactionTable/OldTransactionTable*. *ClosedTable* which is used to maintain the information of closed itemsets consists of three fields: *Cid*, *CP* and *SCP*. Each closed itemset *CP* is assigned to a unique closed identifier *Cid*, and its support count is denoted as *SCP*. *CidList* maintains each item in data stream and its corresponding *Cid* sets which point to *ClosedTable*. *NewTransactionTable* is used to maintain the information of new transaction T_{new} . It consists of two fields: *TempItem* and *Cid*. *TempItem* contains the information of itemsets which satisfy $\{T_i \cap T_{new}, T_i \in \text{ClosdeTable}\}$. The T_i is the i th transaction in data stream and T_{new} is the new transaction. Structure of

OldTransactionTable is same as *NewTransactionTable*, which is used to maintain the information of old transaction T_{old} .

The core issue of removing old transactions from sliding window is how to effectively prune the existing data structures. Existing methods are often pruning step by step which is inefficient. A sliding step M is used in this paper. Pruning data structure after the sliding window move M transactions, that is, pruning when transactions T_{sw+i*M} (sw is the size of sliding window, $i=1, 2, \dots$) are arrived.

In order to distinguish the weights of the historical transactions and the recent transactions, thereby improving the accuracy of result set, and avoiding the missing of possible frequent patterns, a novel algorithm TDMCS (TDM-Based Closed Frequent Pattern Mining on Data Stream) is proposed in this paper. TDMCS mines closed frequent patterns on data stream based on frame θ - ϵ - f (decay factor-minimum support-maximal support error). This algorithm uses data structures *ClosedTable*, *CidList*, *NewTransactionTable* and *OldTransactionTable* to maintain frequent itemsets information. It uses time decay model (TDM) to estimate the support count of pattern, and maintain the frequent and half-frequent closed frequent itemsets which satisfy frame θ - ϵ . The description of algorithm TDMCS is shown as Algorithm 1. The main idea is processing the information of new transaction T_{new} at first. Secondly, if the number of processed transactions exceeds the size of sliding window, delete the information of old transaction T_{old} . If the processing steps of the transactions meet the pruning step M , do pruning operation. In order to increase the efficiency of the algorithm, it only adds delete flags (*DeleteFlag*) when processing old transactions and does the actual delete operations when pruning.

Algorithm 1: TDMCS()

Mining closed frequent patterns on data streams

- 1 For Each Transaction T_{new} In S Do
- 2 Call TDMCSADD(T_{new});
- 3 If $NUM > N$ Then TDMCSREMOVE(T_{old});
- 4 If $NUM \% M = 0$ Then Call PRUNNING();
- 5 End For

Specifically, there are three methods consisted in algorithm TDMCS. Method *TDMCSADD*(T_{new}) is used to process the new transactions, method *TDMCSREMOVE*(T_{old}) processes old transactions and method *PRUNNING*() processes pruning.

When new transaction T_{new} arrived, *TDMCS-ADD*(T_{new}) is described as Algorithm 2. For example, if new transaction is T_4 as shown in Table 1, this algorithm processes it in four steps. First, generate *CidSet* associated with T_4 to discover the intersection between T_4 and existed frequent itemsets. Second, build *NewTransactionTable* to maintain possible frequent itemsets associated with T_4 . Then update

ClosedTable and *CidList* referring to *NewTransactionTable* and *ClosedTable*. The processing is as shown in Figure 1.

Algorithm 2: *TDMCSADD()*

Processing new transactions

```

1 Add  $T_{new}$  To NewTransactionTable
2 Let  $setcid(T_{new}) = \{ \cup CidSet(item_i), item_i \in T_{new} \}$ 
3 For  $C_{id}$  In  $setcid(T_{new})$  Do
3.1  $interS = T_{new} \cap ClosedTable(C_{id})$ 
3.2 For TempItem In NewTransactionTable Do
    If  $interS \in ClosedTable$ 
    Then update  $support(interS)$ 
    Else If  $support(interS) \geq N \times \epsilon \times \theta$ 
    Then Add ( $interS, C_{id}$ ) To ClosedTable
End For
3.3 For (TempItem, Cid) In NewTransactionTable Do
    If (TempItem == ClosedTable(Cid))
    Then update  $support(ClosedTable(C_{id}))$ 
    Else update  $support(TempItem)$ 
    If ( $newsupport(TempItem) \geq N \times \epsilon \times \theta$ )
    Then Add (TempItem, newsupport(TempItem))
    To ClosedTable
    If  $item \in T_{new}$  And item Is Not In CidList
    Then Add item To CidList
End For
4 End For

```

Illustrate the process of algorithm *TDMCSADD*(T_{new}) to handle new transactions. Data stream is shown as in Table 1 including 4 transactions. Let decay factor $f=0.8$, minimum support threshold $\theta=0.1$. When new transaction $T_4=\{2, 3, 4, 5\}$ is arrived, information of *ClosedTable*, *CidList* and *NewTransactionTable* are as shown in Tables 2-4. There are 5 frequent itemsets in *ClosedTable*, 5 items in *CidList* and NULL in initialized *NewTransactionTable*. When new transaction T_4 arrived, there are some steps to process it.

Step 0: Add values $\langle T_4, 0 \rangle$ to *NewTransactionTable*.

Step 1: Compare items in T_4 and items in *CidList* to get the *CidSet* associated with T_4 .

Step 2: Add itemsets associated with T_4 to *NewTransactionTable* according with *CidSet*, as shown in Table 4. That is for each element *Cid* of *CidSet* to get the intersection of T_4 and *ClosedTable*.

Step 3: Update *ClosedTable* with information of *NewTransactionTable*, then get Table 5.

Step 4: Update *CidList* with information of novel *NewTransactionTable* and *ClosedTable*. It will be update under two conditions: the emergence of a new frequent itemset or a new item. Assuming T_4 contains a new item (7) represented in italics in Figure 1. Then add value $\langle \{7\}, \{6\} \rangle$ to *CidList*. Meanwhile, there are two new frequent itemsets in *ClosedTable*, then update *CidList* too.

For continuous generation of new transactions, repeat steps above for processing. When you need to remove the information of old transactions from sliding window, use algorithm *TDMCSREMOVE*(T_{old}). The

main process is similar to algorithm *TDMCSADD*(T_{new}). First, generate *OldTransactionTable* to maintain the information about old transaction T_{old} . Second, find the intersections of *OldTransactionTable* and *ClosedTable*. Next, update or delete *ClosedTable*. In order to improve efficiency of algorithms, it only adds deleting flags and does not do the actual deletion.

When steps of processing transactions meet the pruning step, call function *PRUNING*() to do pruning operations. This is the actual process of deleting operations, and it mainly does updating and deleting operations on *ClosedTable* and *CidList*. The algorithm is described as shown in Algorithm 3.

Algorithm 3: *PRUNING()*

Dropping information of historical transactions.

```

1 For Each  $C_{id}$  In ClosedTable
2 Remove itemsets (with DeleteFlag) From closedTable
3 If  $support(C_{id}) < N \times \epsilon \times \theta$ 
    Then Remove itemsets From closedTable
4 Update CidList
5 End For

```

If only the parameters minimum support threshold θ and decay factor f are used in algorithms, some possible frequent patterns might be lost. Such as, let minimum support $\theta=0.1$, then complete result set is mined. Therefore, many useless patterns may be discovered. If setting $\theta=0.3$, when T_4 arrived, frequent itemsets should meet the support count $4 \times 0.3 = 1.2$. Then generate three frequent itemsets in *ClosedTable* as shown in Table 7. From Table 7 and Table 5, it is clear that the pattern $\{3\ 4\}$ ($freq_d(\{3\ 4\}) = 1.512 > 1.2$) is missing.

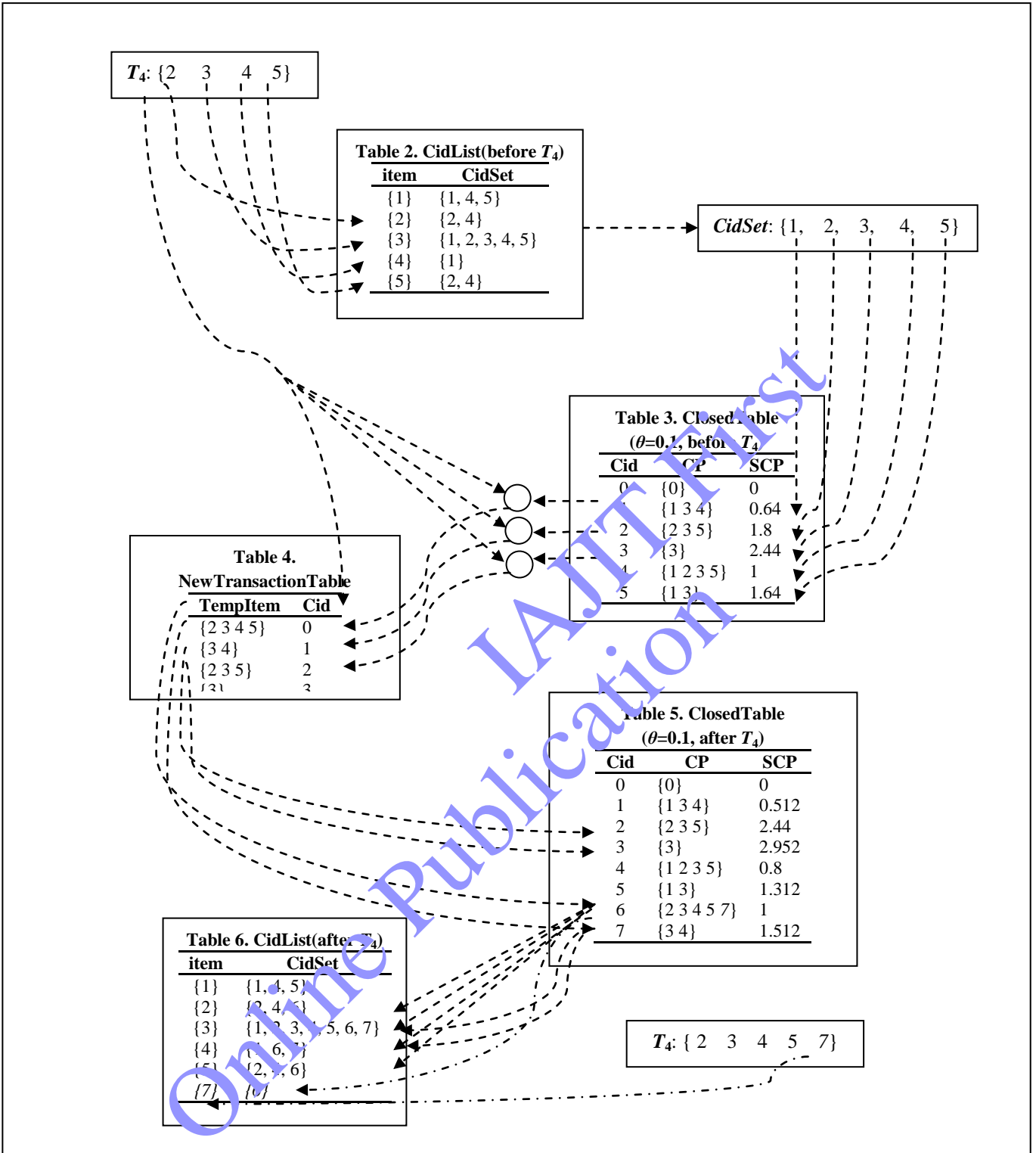


Figure 1. Schematic of the process of handling new transaction T_4

Table 7. ClosedTable ($\theta=0.3$)

Cid	CP	SCP
0	{0}	0
1	{2 3 5}	2.44
2	{3}	2.952
3	{1 3}	1.312

The reason is that the frequency with decay factor of pattern P is smaller than its original frequency, that is $freq_d(P) < freq(P)$. Let $f=0.8$, then $freq_d(P) < 1/(1-0.8)=5$

as calculated by Formula 6. Therefore, if only the minimum support threshold θ is used under the time decay model, some frequent patterns may be lost, for its support count may be less than $\theta \times N$. To solve this problem, ε is introduced as the maximum support error. Therefore, the frequent support of mined patterns needs to meet $N \times \varepsilon \times \theta$ instead of $N \times \theta$.

$$freq_d(P, T_m) = freq_d(P, T_{m-1}) \times f + r \quad (6)$$

$$= \sum_i r_i \times f^{m-i} = r_1 \times f^{m-1} + r_2 \times f^{m-2} + \dots + r_m$$

$$\leq f^{m-1} + f^{m-2} + \dots + 1 \leq \frac{1}{1-f}$$

The next question is how to determine the value of the decay rate f after given parameters: minimum support, maximum support error and sliding window size. Suppose recall is 100%, a lower bound for the decay factor is showed by Formula 7 [3, 12]. Formula 8 [3] shows the upper bound of f under the condition of precision=100%. The usual methods set f to random value between lower bound and high bound or set f to one bound of them[3, 12]. Because both recall=100% and precision=100% cannot be achieved at same time, selected f should balance these two conditions.

$$f \geq \frac{(2N-\theta N^{-1})\sqrt{[(\theta-\varepsilon)/\theta]^2}}{1}, \text{ when recall}=100\% \quad (7)$$

$$f < \frac{(\theta-\varepsilon)N-1}{(\theta-\varepsilon)N}, \text{ when precision}=100\% \quad (8)$$

In this paper, a new way to set decay factor is proposed. Let sliding window size be 10K. Parameters θ , ε and f are shown in Table 8. The third column f_{recall} means the lower bound when assuming recall is 100%. The last column $f_{\text{precision}}$ implies the upper bound when assuming precision is 100%. There are three policies to select f , as shown in Formula 9. For example, let $\theta=0.025$ and $\varepsilon=0.05$, then set $f=f_{\text{recall}}=0.999995$, $f=f_{\text{precision}}=0.995789$ or $f=f_{\text{average}}=0.997892$. Verified by experiments (in Section 4), set f to f_{average} can get the more balanced recall and precision of algorithm. Therefore, set $f=f_{\text{average}}$ is more reasonable than set f on random value between f_{recall} and $f_{\text{precision}}$ or one of it.

Table 8. Time decay factor

θ	$\varepsilon \times \theta$	frecall (recall=100%)	fprecision (precision=100%)
0.05	0.05×0	0.999995	0.997895
0.05	0.1×0	0.999989	0.997778
0.05	0.5×0	0.999929	0.996
0.025	0.05×0	0.999995	0.995789
0.025	0.1×0	0.999989	0.995556
0.025	0.5×0	0.999929	0.992

$$f_1 = f_{\text{recall}} = \frac{(2N-\theta N^{-1})\sqrt{[(\theta-\varepsilon)/\theta]^2}}{1}$$

$$f_2 = f_{\text{precision}} = \frac{(\theta-\varepsilon)N-1}{(\theta-\varepsilon)N} \quad (9)$$

$$f_3 = \frac{(2N-\theta N^{-1})\sqrt{[(\theta-\varepsilon)/\theta]^2} + (\theta-\varepsilon)N-1}{2}$$

4. Performance

The experiments were performed on a 2.1 GHZ CPU with 2GB memory, and running on Win7. All the algorithms were coded in Java language. To evaluate the performance of these algorithms, real and synthetic datasets were used. Real dataset from UCI[7] describes

the page visits of users who visited msnbc.com on September 28, 1999. Visits are recorded at the level of URL category (see below) and are recorded in order. There are 989818 transactions and the average length of transaction is 5.7. It is high dense and similar data stream. Synthetic datasets were generated from IBM data generator. There are four synthetic datasets with different average pattern and transaction size: T5I5D1000K, T10I4D1000K, T10I5D1000K, T10I10D1000K, T20I5D1000K and T20I20D1000K. These were used to analyze the performance of data stream of different density. The parameters are described as follows: D is the total number of transactions; I is the average size of maximal potential patterns; T is the average length of transactions. Such as, T10I5D1000K means the average length of transactions is 10, average length of maximal potential patterns is 5, and number of transactions is 1000K.

The mainly purpose of experimental was to analyze: (1) the ways to set decay factor f . Compared the algorithm performances with setting f as random value, boundary value and average value. (2) Analyzed the effects of sliding window sizes on performance of algorithm TDMCS. (3) Analyzed the effects of pruning steps on performance of algorithm TDMCS. (4) Compared the performances of algorithms TDMCS and CloStream*[25], MSW[12] and SWP[3]. Compared to algorithm CloStream[24], algorithm CloStream* used the sliding window to deal with recent transactions to mine closed frequent patterns. CloStream handled all the transactions, so it did not apply to mine unlimited data stream. Therefore, in this paper compared TDMCS with CloStream*. Similar pattern tree structures were used in algorithms MSW and SWP. And both of them set decay factor as lower bound. In this paper, some modifications were made to the two original algorithms for mining closed frequent patterns instead of complete patterns.

The maximum support error ε was set to 0.1. The value of the time decay factor f was set to average of low bound and high bound to balance 100% recall and 100% precision. The sliding window sizes N were set from 0.1M to 0.8M and the minimum support thresholds were set from 0.06 to 0.1. The values of f in the experiments are shown in Table 9.

Table 9. Values of decay factors

fid	θ	$\varepsilon \times \theta$	N	f
f ₁	0.06	0.006	0.1M	0.990686
f ₂	0.06	0.006	0.2M	0.995343
f ₃	0.06	0.006	0.3M	0.996895
f ₄	0.06	0.006	0.4M	0.997672
f ₅	0.06	0.006	0.5M	0.998137
f ₆	0.06	0.006	0.7M	0.998669
f ₇	0.06	0.006	0.8M	0.998836
f ₈	0.07	0.007	0.1M	0.992009
f ₉	0.08	0.008	0.1M	0.993001
f ₁₀	0.09	0.009	0.1M	0.993772
f ₁₁	0.1	0.01	0.1M	0.994389

At first, verify the reasonableness of setting $f=f_{\text{average}}$. The relationship between decay factor f and minimum support threshold θ , maximal support error threshold ε needs to be discussed in algorithm TDMCS, in order to determine the optimum parameter value of f .

Let minimum support $\theta=0.05$, the values of recall and precision of TDMCS on *msnbc* with different decay factors are shown in Figure 2. Abscissa axis means the random value between f_{recall} and $f_{\text{precision}}$ at different window size N . Vertical axis means the recall and precision at different N , and the dashed line means to set $f=f_{\text{average}}$. From this figure, it can be concluded that: (1) with the decreasing of f , recall is decreasing and precision is increasing. (2) The trends of recall and precision at different N are similar. (3) When setting $f=f_{\text{average}}$, the values of recall and precision are fixed. They are almost unaffected by sizes of sliding windows. (4) The values of recall and precision can be balanced by setting $f=f_{\text{average}}$.

When setting $f=f_{\text{recall}}$, $f=f_{\text{precision}}$ and $f=f_{\text{average}}$, compare the average values of recalls and precisions of algorithm with different sliding window sizes. It can be concluded that it can get almost 100% recall when setting $f=f_{\text{recall}}$, and get lowest recall with setting $f=f_{\text{precision}}$. When setting $f=f_{\text{average}}$, the value of recall is between them. Setting $f=f_{\text{precision}}$ and $f=f_{\text{average}}$ can get almost the same precisions. But the value of precision is lowest when setting $f=f_{\text{recall}}$. Therefore, recall and precision of algorithm can be more balanced by setting $f=f_{\text{average}}$ than setting $f=f_{\text{recall}}$ and $f=f_{\text{precision}}$.

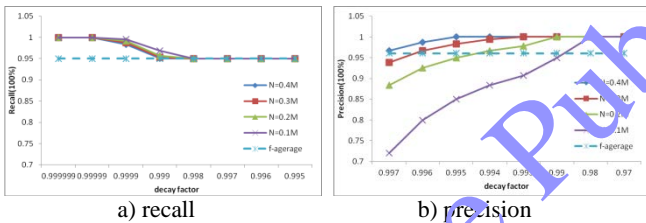


Figure 2. Variation of recall and precision with different decay factor

Next, compare the performance of algorithm with setting f to f_{average} and random values. To make the random value more reasonable, set them in the range of (0.9, 1), and denoted as f_{random} . Use function *Math.random()* to generate 5 random values to set decay factors. The performance of TDMCS on *msnbc* is shown in Figure 3. As can be seen, when setting $f=f_{\text{average}}$ and $f=f_{\text{random}}$, the values of precision are little different. But the values of recall are very different when setting f to random values. Therefore the performance of algorithm is unstable. The performance with $f=f_{\text{average}}$ is significantly better than set $f=f_{\text{random}}$, and the result set obtained is stable.

It can get the same conclusions when processing synthetic data streams. Thus, set decay factor to average value is reasonable.

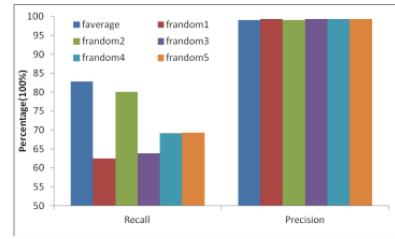


Figure 3. The performance studies on *msnbc* with set decay factor to *faverage* and random values

The second performance comparison experiments are used to analyze the influence of window sizes on algorithm TDMCS. Set sliding window size $N=0.1M$, $0.2M$ and $0.3M$; the minimum support $\theta=0.06$; the decay factor $f=f_1, f_2, f_3$ as shown in Table 9; pruning step $P=0.1M$ [3, 12]. The runtime and space costs of algorithm TDMCS of *msnbc* are compared in Figure 4.

The performance of TDMCS on data stream *msnbc* is shown in Figure 4 when processing 1M, 1.5M, 2M and 2.5M transactions. Figure 4.a shows the runtime and in which abscissa axis means number of transactions. It can be seen that: (1) when the number of transactions is small, the increase in window size leads to a slight increase in runtime; (2) with the increase of processing transaction number, the runtime with bigger window sizes is lower than runtime of smaller window size. Figure 4.b shows the memory usage. It is clear that the effect of different window size on memory usage is small. From time and space consumption, it can be concluded that the runtime of algorithm TDMCS on data stream *msnbc* is greatly different as different size N under the same number of transactions. And the memory usage is almost same as different size N . Thence, in terms of space complexity, TDMCS applies to discovering frequent patterns of any window size.

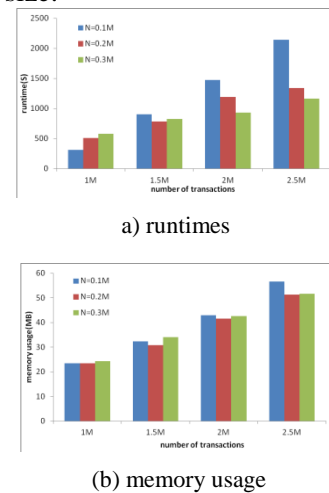


Figure 4. The performance study of TDMCS on *msnbc* based on different sizes of windows

Thirdly, analyze the effect of pruning step on performance of algorithm TDMCS on *msnbc*. Sliding window size N was set to $0.5M$, $0.7M$ or $0.8M$. Pruning step P was set to $0.1M$ to $0.5M$ ($P \leq N$).

Minimum support θ was set to 0.06 and f was set to f_5 to f_7 .

Figure 5.a shows the runtime of TDMCS performance on data stream *msnbc* with different window sizes and different pruning step lengths. From the performance of runtime it can be concluded that: (1) optimal pruning step length is related to sliding window size, (2) when value P is different, the runtime consumption is obvious different as N increasing. The memory usage is shown in Figure 5.b. It is clear that the effect of pruning step length on memory usage is small. From Figure 5 it can be seen that the length of pruning step has little influence on runtime and memory usage when window size is small. And when window size is big, pruning step has a wider range influence on runtime. Therefore, set the parameter $P=0.1M$ when $N=0.1M$ is reasonable. Figure 5.b also shows that algorithm TDMCS applies to discovering frequent patterns of any size of window.

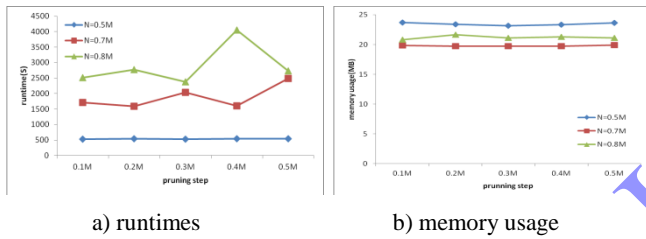


Figure 5. The performance study of TDMCS on data stream *msnbc* based on different sizes of windows and pruning steps.

Finally, compare the performances of TDMCS and classical algorithms with different sliding window sizes. Let minimum support $\theta=0.06$, maximal support error $\varepsilon=0.1$, pruning step $P=0.1M$, sliding window size $N=0.1M$ to $0.5M$ and decay factor $f=f_1$ to f_7 as shown in Table 9.

The performances of four algorithms processing synthetic data streams are shown in Figure 6, which are average values of recalls and precision under different sliding window sizes. Four data streams with different lengths of transactions or patterns are used, including: T10I4, T10I5, T10I10 and T20I5. Figure 6.a shows the runtimes of four algorithms. Overall, the runtimes of TDMCS and CloStream* are lower than other two algorithms. This is because algorithm CloStream* does not process data with decay operations and algorithm TDMCS uses closure operator. The memory usages of algorithms are shown in Figure 6.b. The memory usage of TDMCS is the lowest of all. But the different between four algorithms is not too much. Figure 6.c and Figure 6.d show the recalls and precisions of algorithms. Both algorithms MSW and SWP set decay factor to lower bound, so we only compared with SWP. The recall of algorithm CloStream* is highest of all for it does not use time decay model, but the precision is the lowest of all. The recall and precision of algorithm SWP are in the middle of three algorithms. The recall of TDMCS is about 1% lower than other two

algorithms, but the precision of TDMCS is about 10% higher than CloStream* and about 4% higher than SWP. Therefore, algorithm TDMCS can get more balance recall and precision than other three methods. And compared performances of four algorithms on data streams with different pattern lengths, such as T10I4, T10I5 and T10I10, or data streams with different transaction lengths, such as T10I5 and T20I5, it can be concluded that algorithm TDMCS is more suitable to process data streams with long transaction length and long pattern length.

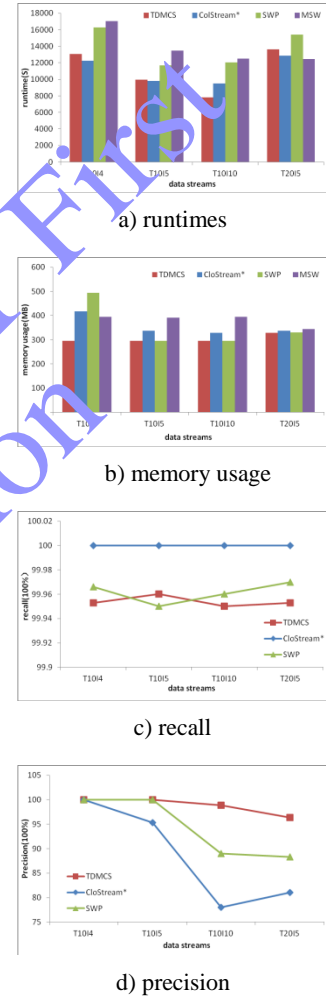


Figure 6. The performance studies of algorithms on synthetic data streams.

5. Summary

Data stream is a fluid, continuous, unbounded and time ordered sequence of data transactions generated at a rapid rate. Due to the knowledge contained in data stream will change over time, concept drift should be taken into account when mining frequent patterns. Normally, recent transactions contain more important information than historical transactions, thus they should be treated differently. Considering the data stream characteristics, an efficient algorithm TDMCS is proposed in this paper. It is used to mining closed frequent patterns and based on time decay model and

sliding window model. It uses closure operator to improve the efficiency of mining closed frequent itemsets. In order to balance the recall and precision, a novel manner by average the low bound and high bound is provided in this paper. It uses frame minimum support-maximum support error-decay factor to avoid concept drift and discover more reasonable and compact result set. The performance of the proposed algorithm was investigated using experiments. The results show that it is efficient and scalable, and it applies to mining high dense data stream and long patterns.

References

- [1] Chang T. P. "Mining frequent user query patterns from xml query streams," *International Arab Journal of Information Technology*, Vol. 11, No. 5, 2014: 452-458.
- [2] Chen H. "Mining top-k frequent patterns over data streams sliding window," *Journal of Intelligence Information System*, Vol. 42, Issue 1, 2014: 111-131.
- [3] Chen H, Shu L C, Xia J L, and Deng Q S. "Mining frequent patterns in a varying-size sliding window of online transactional data streams," *Information Sciences*, 2012, 215(12): 15-36.
- [4] Cheng J, Ke Y, and Ng W. "Maintaining frequent closed itemsets over a sliding window," *Journal of Intelligent Information Systems*, 2008, 31(3): 191-215.
- [5] Chi Y, Wang H X, Yu P S, and Muntz K R. "Catch the moment: maintaining closed frequent itemsets over a data stream sliding window," *Knowledge and Information Systems*, 2006, 10(3): 265-294.
- [6] Farzanyar Z., Kangavari M., and Cercone N.. "Max-FISM: Mining (recently) maximal frequent itemsets over data streams using the sliding window model," *Computers and Mathematics with Applications*, Vol. 64, 2012: 1706-1718.
- [7] Frank A, and Asuncion A. UCI Machine Learning Repository[EB/OL]. Irvine, CA: University of California, School of Information and Computer Science, <http://archive.ics.uci.edu/ml>, 2010.
- [8] Hewa Nadungodage C., Xia Y., Lee J. J., and Tu Y.. "Hyper-structure mining of frequent patterns in uncertain data streams," *Knowledge and Information Systems*, Volume 37, Issue 1, 2013: 219-244.
- [9] Jiang N and Gruenwald L. "CFI-Stream: mining closed frequent itemsets in data streams," *Proceedings of ACM SIGKDD Internal Conference on Knowledge Discovering and Data Mining*, New York, USA, 2006: 592-597.
- [10] Lee G. , Yun U. , and Ryu K. H.. "Sliding window based weighted maximal frequent pattern mining over data streams," *Expert Systems with Applications*, Volume 41, 2014: 694-708.
- [11] Li H. F., Ho C. C., and Lee S. Y.. "Incremental updates of closed frequent itemsets over continuous data streams," *Expert Systems with Applications*, Vol. 36, Issue 2, 2009: 2451-2458.
- [12] Li G. H., and Chen H. "Mining the frequent patterns in an arbitrary sliding window over online data streams," *Journal of Software*, 2008, 19(19): 2585-2596.
- [13] Li H. F., Zhang N, et al. "Frequent itemset mining over time-sensitive streams," *Chinese Journal of Computers*, Vol. 35, No. 11, 2012: 2283-2293.
- [14] Li H. F., Ho C. C., Chen H. S., and Lee S. Y.. "A single-scan algorithm for mining sequential patterns from data streams," *International Journal of Innovative Computing, Information and Control*, Volume 8, Number 3(A), 2012: 1799-1820.
- [15] Manu Q, and Motwani. "Approximate frequency counts over streaming data," *Proceedings of the 28th International Conference on Very Large Data Bases*, Hong Kong, China, 2002: 346-357.
- [16] Nabil H. M., Eldin A. S., and Belal M. A. E.. "Mining frequent itemsets from online data streams: comparative study," *International Journal of Advanced Computer Science and Applications*, Vol. 4, No.7, 2013: 117-125.
- [17] Nori F, Deypir M, and Sadreddini M H. "A sliding window based algorithm for frequent closed itemset mining over data streams," *Journal of Systems and Software*, 2013, 86(3): 615-623.
- [18] Patnaik D., Laxman S., Chandramouli B., and Ramakrishnan N.. "A general streaming algorithm for pattern discovery," *Knowledge and Information Systems*, Vol. 37, Issue 3, 2013: 585-610.
- [19] Shie, B. E., Yu, P. S., and Tseng, V. S.. "Efficient algorithms for mining maximal high utility itemsets from data streams with different models," *Expert Systems with Applications*, Vol. 39, 2012: 12947-12960.
- [20] Tang K M , Dai C Y, and Chen L. "A novel strategy for mining frequent closed itemsets in data streams," *Journal of Computers*, 2012, 7(7): 1564-1572.
- [21] Tsai P. S. M.. "Mining top-k frequent closed itemsets over data streams using the sliding window model," *Expert Systems with Applications*, Vol. 37, Issue 10, 2010: 6968-6973.
- [22] Wong R. C. W., and Fu A. W. C.. "Mining top-k frequent itemsets from data streams," *Data Mining and Knowledge Discovery*, Vol.13, Issue 2, 2006: 193-217.

- [23] Yang B., and Huang H.. "TOPSIL-Miner: an efficient algorithm for mining top-k significant itemsets over data streams," *Knowledge and Information Systems*, Vol. 23, Issue 2, 2010: 225-242.
- [24] Yen S J, Lee Y S, Wu C W, and Lin C L. "An efficient algorithm for maintaining frequent closed itemsets over data stream," *Next-Generation Applied Intelligence*, 2009, 5579(1): 767-776.
- [25] Yen S J, Wu C W, and Lee Y S et al. "A fast algorithm for mining frequent closed itemsets over stream sliding window," *Proceedings of 2011 IEEE International Conference on Fuzzy Systems*, Taipei, Taiwan, 2011: 996-1002.
- [26] Yu J X, Chong Z, Lu H, and Zhou A. "False positive or false negative: mining frequent itemsets from high speed transactional data streams," *Proceedings of the Thirtieth International Conference on Very Large Data Bases*, Toronto, Canada, 2004: 204-215.



Han Meng, born in 1982, Ph.D. candidate, associate professor. Her research interests include data mining and machine learning.



Jian Ding, born in 1977, M.S., associate professor. His research interests include machine learning and data mining.



Juan Li, born in 1975, M.S., associate professor. Her research interests include information security and cloud computing.

Online Publication
IAJIT First