

会话流中 Top-k 闭序列模式的挖掘

彭慧丽¹, 张啸剑²

(1. 河南省直广播电视大学教务科, 郑州 450008; 2. 河南财经学院计算机系, 郑州 450002)

摘 要: 在会话流中挖掘 Top-k 闭序列模式, 存在因相关比率 ρ 的大小而导致的内存消耗和挖掘精度之间的冲突。基于 False-Negative 方法, 提出 Tstream 算法, 制定 2 种约束策略限制 ρ 。基于该策略设计加权调和计数函数, 渐进计算每个模式的支持度。实验结果证明了该算法的有效性。

关键词: Top-k 闭序列模式; 加权调和平均数; 调节因子

Top-k Closed Sequential Pattern Mining in Session Streams

PENG Hui-li¹, ZHANG Xiao-jian²

(1. Department of Education, Henan Radio & Television University, Zhengzhou 450008;

2. Department of Computer Science, Henan University of Finance & Economics, Zhengzhou 450002)

【Abstract】 The current methods in session streams for mining Top-k Closed Sequential Pattern(Topk_CSP) may lead to a conflict between output precision and memory consumption because of using ρ . This paper proposes TStream algorithm, which is based on False-Negative approach. TStream utilizes two constraint strategies to restrict ρ , and employs a weighted harmonic count function to calculate the support of each pattern progressively. Experimental results show that the algorithm is efficient.

【Key words】 Top-k Closed Sequential Pattern(Topk_CSP); Weighted Harmonic Average(WHA); regulatory factor

1 概述

Top-k 闭模式挖掘一直是数据流研究中的热点。False-Positive 和 False-Negative 是常用的 2 类模式挖掘方法。由于会话流具有数据流的连续性、无界性等特点, 传统的 Top-k 闭序列模式(Top-k Closed Sequential Pattern, Topk_CSP)挖掘算法已不适用。研究者们针对数据流提出了许多算法挖掘频繁模式^[1-2]和 Top-k 频繁模式^[3]。但这些算法存在如下须解决的问题: (1)利用类似文献[1-2]中的算法通常会产生大量的模式。(2)采用文献[2]中方法挖掘闭模式时, 支持度阈值 σ 的设置非常敏感。 σ 值太小会导致过多的闭模式; σ 值太大会导致无闭模式产生。(3)目前大多数算法基于 False-Positive 方法挖掘各种模式, 如文献[1, 3]中的算法。False-Positive 方法利用相关比率 ρ 控制内存消耗、挖掘精度和查全率。使用较大的 ρ 会降低内存消耗, 但精确性降低; 使用较小的 ρ 能够提高精度, 但内存消耗增加, 挖掘效率降低。

为解决上述问题, 本文挖掘会话流中的 Topk_CSP, 制定 2 个边界参数来限制 ρ , 设计 2 个边界的加权调和平均数替代 ρ , 并设置 1 个调节因子调节 ρ 值的大小。在此基础上提出一种基于 False-Negative 方法和时间敏感滑动窗的挖掘算法 TStream, 有效地挖掘某一会话流上的 Topk_CSP。

2 相关概念和描述

令 $P=\{P_1, P_2, \dots, P_n\}$ 是 Web 页面的完全集合。一个会话 S 是一个由时间戳指定顺序的序列。一个序列是由被访问时间标记的 Web 页面组成的。会话流 S_s 是由不断到达的会话组成的动态增长会话集, 即 $S_s=\{S_1, S_2, \dots, S_m, \dots\}$ 。时间敏感滑动窗 Tsw 是一个向前滑动的窗口, 由一组连续的时间单元组成的集合。设当前滑动窗为 Tsw_T , 则 $Tsw_T=\langle t_{T-w+1}, t_{T-w+2}, \dots,$

$t_T \rangle$, 其中, t_T 为当前时间单元; 窗口的大小为 w 。在窗口 Tsw 中, 当序列 s 满足条件 $C(s, Tsw) \geq \sigma |Sset(Tsw)|$ 时, s 为频繁序列模式, 其中, $Sset(Tsw)$ 表示在 Tsw 中到达的会话集合; σ 为给定的一个支持度阈值; $C(s, Tsw)$ 表示 $Sset(Tsw)$ 中包序列 s 的会话数目。

定义 1 给定 Tsw 和序列模式 s , 如果不存在这样的模式 s' , 使得条件 $s \subset s'$ 和 $C(s, Tsw) = C(s', Tsw)$ 同时满足, 则 s 为闭序列模式。如果恰好存在 $(k-1)$ 个闭序列模式 $s_i'' (i=1, 2, \dots, k-1)$, 满足 $C(s_i'', Tsw) > C(s, Tsw)$, 则 s 在 Tsw 上是 Topk_CSP。

3 约束方法和加权调和计数函数

由于会话流的自身特性, 挖掘其中的 Topk_CSP 会出现一定的误差, 主要分为面对挖掘精度和面对查全率 2 种。目前, 许多基于 False-Positive 方法的算法采用 ρ 来控制这 2 个误差。然而, 使用 ρ 会在挖掘精度、内存消耗和查全率之间产生矛盾。

3.1 约束方法

为了解决上述矛盾, 给出 2 个边界参数 λ_1 和 λ_2 来约束 ρ , 并且制定约束策略来满足用户的挖掘目的。

(1)如果 $\rho < \lambda_1$, 则触发第 2 种误差。一个很小的 ρ 值能生成大量的候选模式, 导致内存消耗增加, 挖掘效率降低, 因此, $\rho > \lambda_1$ 。

(2)如果 $\rho > \lambda_2$, 则触发第 1 种误差。由于很大的 ρ 将会

基金项目: 河南省科技厅基金资助项目“非线性降维技术在商业智能中的应用”(082300410110)

作者简介: 彭慧丽(1981-), 女, 硕士, 主研方向: 数据挖掘, 数据模型; 张啸剑, 硕士

收稿日期: 2009-01-25 **E-mail:** xjzhang82@yahoo.com

产生精度很低的输出结果, 因此 $\rho < \lambda_2$ 。

3.2 加权调和计数函数

最大化($\rho \approx \lambda_1$)或最小化($\rho \approx \lambda_2$)均会加剧上述 2 种误差, 因此, 设计参数 λ_1 和 λ_2 的加权调和平均数(WHA)来代替 ρ , 即 $\rho = WHA(\lambda_1, \lambda_2)$,

$$WHA(\lambda_1, \lambda_2) = (1 + \xi^2) \lambda_1 \lambda_2 / (\lambda_1 + \xi^2 \lambda_2) \quad (1)$$

而在大多数基于 False-Positive 方法的算法中, ρ 等于 ϵ/σ , 其中, ϵ 为误差参数。采用 $WHA(\lambda_1, \lambda_2)$ 替代 ρ , 则等式 $\rho = \epsilon/\sigma$ 将发生变化,

$$WHA(\lambda_1, \lambda_2) = \epsilon/\sigma \quad (2)$$

$$\epsilon = \sigma \times (1 + \xi^2) \lambda_1 \lambda_2 / (\lambda_1 + \xi^2 \lambda_2) \quad (3)$$

式(1)和式(3)中的参数 ξ 是一个调节因子, 通过调整参数 ξ 的值来调节上述 2 种误差以及克服 ρ 引起的问题。

定义 2 根据式(3), 序列模式 s 在一个时间单元 t_i 上的潜在支持度计数定义如下:

$$\hat{C}(s, t_i) = \begin{cases} 0 & C(s, t_i) < \epsilon | Sset(t_i) | \\ C(s, t_i) & \text{其他} \end{cases} \quad (4)$$

因此, 序列模式 s 在当前滑动窗 $Tsw_I = \langle t_{I-w+1}, t_{I-w+2}, \dots, t_I \rangle$ 上的累积支持度计数定义如下:

$$\hat{C}(s, Tsw_I) = \sum \hat{C}(s, t_i) \quad (5)$$

其中, $t_i \in Tsw_I, I-w+1 \leq i \leq I$ 。

定义 3 给定参数 λ_1 和 λ_2 , $Tsw_I = \langle t_{I-w+1}, t_{I-w+2}, \dots, t_I \rangle$ 为当前的滑动窗。令 $\langle t_{I-R+1}, t_{I-R+2}, \dots, t_I \rangle$ 为当前窗口 Tsw_I 中最近出现的 R 个时间单元, 命名为 Tsw_R , 大小为 $|Sset(Tsw_R)|$, $1 \leq R \leq w$ 。则加权调和计数函数(WHC)定义如下:

$$WHC(R) = \lceil \sigma | Sset(Tsw_R) | \times WHA(\lambda_1, \lambda_2) \rceil \quad (6)$$

可知, 在 Tsw_R 中的序列模式 s , 如果 $\hat{C}(s, Tsw_R) \geq WHC(R)$, 则 s 在 Tsw_I 中为潜在频繁序列模式。否则 s 应从 Tsw_I 删除。

定义 4 给定 Tsw_R 和潜在序列模式 s , 如果不存在这样的模式 s' , 使得 $s \subset s'$ 和 $\hat{C}(s, Tsw_R) = \hat{C}(s', Tsw_R)$ 条件同时成立, 则 s 为潜在闭序列模式。如果恰好存在 $(k-1)$ 个潜在闭序列模式 $s_i'' (i=1, 2, \dots, k-1)$, 满足 $\hat{C}(s_i'', Tsw_R) > \hat{C}(s, Tsw_R)$, 则潜在闭序列模式 s 在 Tsw_R 上为潜在 Top- k 闭序列模式(PTk_CSP)。

4 TStream 算法

TStream 算法包括 2 个子程序: (1)Ttree 构造 Top- k 树; (2)Mtree 对 Top- k 树进行维护。在子程序 Ttree 中函数 *leftcheck* 检测某模式是否为 PTk_CSP。

Top- k 树是一种字典序列树, 类似于前缀树。由 3 个部分组成: (1)Top- k 树由 2 部分组成: 一个带有标记为 \emptyset 的根节点和一个页面前缀子树集合。(2)Top- k 树中除根节点之外的每个节点由 3 个域组成: *page*, *tid*, $\hat{C}(s, Tsw)$, 其中, *page* 记录模式 s 中的最后一个页面; *tid* 记录 t_i 的 *id*, 在 *id* 时刻 s 被插入到树中; $\hat{C}(s, Tsw)$ 表示 s 在窗口 Tsw 中支持度。(3)hash 表用来检测一个模式 s 是否是 PTk_CSP。使用 PTk_CSP 的 $WHC(R)$ 作为哈希地址。

设 $Fset$, $Cset$ 分别是当前时间单元 t_I 和当前窗口 Tsw_I 上的 PTk_CSP 集合。 k 为一个整数。TStream 算法代码如下:

```
Subroutine 1 Ttree (Ss,  $\sigma$ ,  $k$ ,  $\lambda_1$ ,  $\lambda_2$ )
Create root of Top-k Tree T;
foreach  $t_i, t_i \in Tsw_{first} = \langle t_1, t_2, \dots, t_j \rangle$  do
mine all PTk_CSP from Sset( $t_i$ );
if leftcheck ( $s_i, s_i \in PTk\_CSP$ ) = false then
foreach sibling  $s_m$  of  $s_i$  do
create a new child of form ( $s_i \cup m, i, 1$ );
```

```
foreach child  $s_i'$  of  $s_i$  do
Ttree ( $s_i', \sigma, k, \lambda_1, \lambda_2$ );
if ( $s_i \not\subset T$ ) and ( $\hat{C}(s_i, t_i) = \hat{C}(s_i', t_i)$ ) and ( $Fset(t_i, WHC(R)) > k$ )
then
Fset = Fset  $\cup \{s_i\}$ ;
if ( $s_i \subset T$ ) then
add  $\hat{C}(s_i, t_i)$  to  $\hat{C}(s_i)$ ;
if ( $\hat{C}(s_i) < WHC(i-tid(s_i)+1)$ ) or ( $\hat{C}(s_i, Tsw_R) = \hat{C}(s_i', Tsw_R)$ ,
 $s_i \subset s_i'$ ) or ( $Fset(t_i, WHC(R)) < k$ ) then
delete  $s_i$  from T;
Call Mtree (T,  $s_i, \sigma, k, \lambda_1, \lambda_2, w$ );
Subroutine 2 Mtree (T, Ss,  $\sigma, k, \lambda_1, \lambda_2, w$ )
foreach incoming  $t_r (t_r \in Tsw_r)$  do
mine all PTk_CSP from Sset( $t_r$ );
if leftcheck ( $s_i, s_i \in Topk\_CSP$ ) = false then
create a new child of form ( $s_i \cup m, r, 1$ );
foreach child  $s_i'$  of  $s_i$  do
Mtree (T,  $s_i', \sigma, k, \lambda_1, \lambda_2, w$ );
if ( $s_i \not\subset T$ ) and ( $\hat{C}(s_i, t_r) = \hat{C}(s_i', t_r)$ ) and ( $Fset(t_r, WHC(R)) > k$ ) then
update  $\hat{C}(s_i)$  of  $s_i$ ;
Fset = Fset  $\cup \{s_i\}$ ;
if (( $r-tid(s_i)+1 < w$ ) and ( $\hat{C}(s_i) < WHC(r-tid(s_i)+1)$ ) or (( $r-tid(s_i)+1 \geq w$ ) and ( $\hat{C}(s_i) < WHC(w)$ )) or ( $Fset(t_r, WHC(R)) < k$ ) then
delete  $s_i$  from T;
if  $\hat{C}(s_i, Tsw_r) \geq \sigma | Sset(Tsw_r) |$  then
Cset = Cset  $\cup \{s_i\}$ ;
foreach expiring  $t_{r-w+1} (t_{r-w+1} \in Tsw_L)$ 
mine all PTk_CSP from Sset( $t_{r-w+1}$ );
foreach child  $s_i'$  of  $s_i$  do
Mtree ( $s_i'$ );
same as lines 7-13 in Subroutine 2;
if  $s_i \in L_{set}$  and ( $r-tid(s_i)+1 \geq w$ ) then
eliminate  $s_i$  from T;
Output Topk_CSP on demands;
```

5 算法性能分析

TStream 算法采用 C 语言编写, GCC 编译。机器配置是 2.8 GHz Pentium 处理器, 1.0 GB 内存及 Fedora Core 6.0。采用 BMS-WebView-1 和 BMS-WebView-2 数据集。BMS-WebView-1 由 59 602 个会话组成, 平均会话大小为 7 个~13 个页面。BMS-WebView-2 由 537 083 个会话组成, 平均会话大小为 10 个页面。每个 Tsw_i 包含 20 个时间单元, 每个时间单元 t_i 包括近似 100 k 个会话。实验从内存消耗、挖掘精度和查全率 3 个方面测试 TStream 性能。

图 1 显示了在数据集 BMS-Web View-1 上, TStream 算法和基于 Lossy Counting 算法的 Top- k 闭序列模式挖掘算法(TKCLC)在不同 ξ 值下的内存消耗比较。

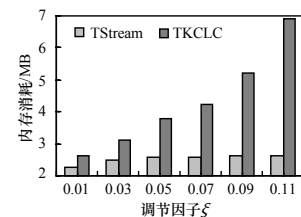


图 1 调节因子变化时的内存消耗

设 $\lambda_1=0.001$, $\lambda_2=0.999$, $\sigma=0.010$, $k=100$ 。结果表明, ξ 从 0.01 变化到 0.11 时, TKCLC 算法的内存消耗增加较快, 而

这些因素与信息处理系统本身紧密结合。因此,该算法在调度过程中,能充分满足系统提出的时限、容量等任务要求,也能适时、适量、适度地将任务推荐给操作员,使调度内容和过程尽可能符合操作员的操作习惯和处理能力。

4.2 实验结果

本算法在分布式测井数据自动处理系统中得到了应用,取得了良好的效果。

验证环境包括 6 台 PC 服务器,配置为:CPU 为双核 2.3 GHz,内存 4 GB,采用 100 Mb/s 以太网互联;工作站/操作员配置为:12 台/人(2 组),CPU 为 1.8 GHz,内存 512 MB,采用 100 Mb/s 以太网互联;为了充分验证本算法的适用程度,选择的操员工龄为 6 个月~14 年不等(职称/熟练程度:实习至助理工程师)。测试数据包括 1.45 GB 的测井数据(通常,一条测井曲线的数据量为 8 MB~10 MB,本次试验的数据被分为 10 组/批次)和其他辅助数据 0.32 GB(主要包括分层数据、曲线修正数据等)。系统先不运行算法,随机分配任务给操作员;10 个工作日后,再运行本算法,由另一组对等的操作员重新处理相同数据,最终得出对比验证结果。

图 3 显示的是在 10 个工作日内,2 种算法处理的测井数据曲线的条数。

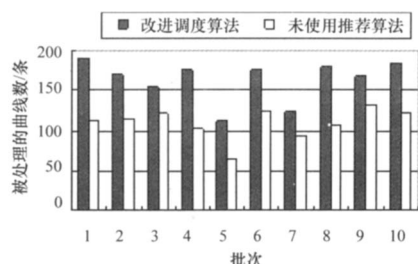


图3 算法性能实验结果

可以看出,应用本算法的系统处理速度明显比未采用本算法的系统快,这主要是由于基于能力感知机制的人际任务

调度算法主动将适宜的、匹配的数据分配给操作员,调动了操作员的能动性和积极性;使得分配的任务和操作员的实际情况相匹配,操作员满意度较高,操作员能够对推荐的任务积极响应;同时使系统和操作员的工作量饱满,缩短了单位任务的响应时间。

5 结束语

本文提出一种基于能力感知的高通量人机任务调度算法,阐述了其主要思想、模型及数据结构。实验表明该算法具有较高的系统利用率和响应率。算法稍作调整,也可应用于其他人机交互频繁的系统,例如某高校的“自动阅卷/评估系统”等。本算法作为对人机交互调度的初期研究成果,目前仅将处理能力和兴趣引入模型,今后可进一步引入其他特征。

参考文献

- [1] Schafer J B, Konstan J, Riedl J. Electronic Commerce Recommender Applications[J]. Journal of Data Mining and Knowledge Discovery, 2001, 5(1/2): 115-152.
- [2] Radulescu A, Van G A. Low-cost Task Scheduling for Distributed-memory Machines[J]. IEEE Trans. on Parallel and Distributed Systems, 2002, 13(6): 648-658.
- [3] 周双娥,袁由光,熊兵周,等. 基于任务复制的处理器预分配算法[J]. 计算机学报, 2004, 27(2): 216-223.
- [4] 黎星星,黄小琴,朱庆生. 电子商务推荐系统研究[J]. 计算机工程与科学, 2004, 26(5): 7-10.
- [5] 王实,高文,李锦涛. 基于分类方法的Web站点实时个性化推荐[J]. 计算机学报, 2002, 25(8): 845-852.
- [6] Bajaj R, Agrawal D P. Improving Scheduling of Tasks in a Heterogeneous Environment[J]. IEEE Trans. on Parallel and Distributed Systems, 2005, 15(2): 107-118.

编辑 顾姣健

(上接第 87 页)

TStream 算法在 ζ 变化时,内存消耗保持稳定。

图 2 和图 3 显示了 TStream 算法和 TKCLC 算法在数据集 BMS-WebView-2 上挖掘精度和查全率的比较。

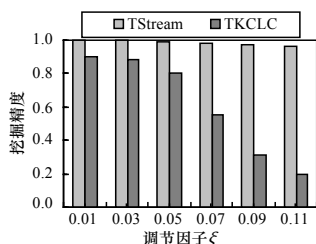


图2 调节因子变化时的挖掘精度

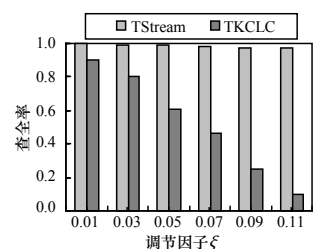


图3 调节因子变化时的查全率

参数 $\lambda_1, \lambda_2, \sigma$ 和 k 与前组实验保持一致。结果表明, ζ 从

0.01 增加到 0.11 时, TStream 算法的平均挖掘精度和查全率几乎接近 97%, 而 TKCLC 算法的挖掘精度和查全率明显降低。

6 结束语

本文提出 TStream 算法挖掘滑动窗口中的 Top_k_CSP, 设计了 Top- k 树结构, 维护窗口中的 Top_k_CSP。在 TStream 算法中, 使用带约束的加权调和平均数处理相关比率引起的问题。当 2 个边界固定以后, 通过调整调节因子的值, 可解决挖掘精度、查全率和内存消耗之间的冲突。

参考文献

- [1] Manku G S, Motwani R. Approximate Frequency Counts over Data Streams[C]//Proc. of VLDB'02. Hong Kong, China: [s. n.], 2002.
- [2] Chi Yun, Wang Haixun, Yu P S, et al. Moment: Maintaining Closed Frequent Itemsets over a Stream Sliding Window[C]//Proc. of ICDM'04. Brighton, UK: [s. n.], 2004.
- [3] Wong R C W, Fu A W C. Mining Top-k Frequent Itemsets from Data Streams[J]. Data Mining and Knowledge Discovery, 2007, 10(13): 193-217.

编辑 顾姣健