

ESWA-FT-
PatternGrowth

Experiments

此处作者使用了4个数据集Retail, BMSWebViewl, FoodMart,T10I4D100k,compare with FT-Apriori, VB-FT-Mine,FT-TreeBased

对于retail和T10I4D100k, 则设置δ=1,2,3,4,5

图9显示4个算法在数据集retail上的执行

图10显示4个算法在数据集BMSWebViewl上的执行

图11显示了4个算法在数据集T10I4D100k上的执行

图12显示了4个算法在数据集FoodMart上的执行

FT-patternGrowth消耗内存还是较高的

图9,10,11,12的C图显示内存消耗情况

scalability analysis

关于数据集每次交易长度

每笔交易数据量大小

第一次计算每个项目的sup,大于小于阈值的会被删除,每条交易中被保留的项目则会根据支持频率重新排序

每一交易都会被映射到FP-tree的分支,公共部分会被合并,两次扫描数据库

第二次扫描数据库,会构建FP-tree树,逐条插入,有相同前缀的会被插入到同一条分支中

项集X是一个频繁项集,算法会产生a,b,c的conditional patterns,算法会忽略其他的conditional patterns,这类似于构建FP-tree的一个分支不包含X中的任何项,那么该分支的FT因子则变为(δ+3),不符合(δ+2),算法会以频繁项集的支持频率率X中的项,对于项集X,FP-tree会生成如下conditional patterns集合

$$\begin{cases} b : < cfab : 1 >, < efc b : 1 >, < eac b : 1 > \\ c : < efc : 1 >, < eac : 1 >, < fdc : 1 > \end{cases}$$

依次找出b,c的条件模式树,若条件模式CB是之前已经发现的条件模式的子集,则对CB的支持度从CA中减去,若CB不为空,则CB被忽略,c中efc和ac为子集被减去,则c中条件模式只剩下dc

$$a : < efa : 2 >, < ea : 1 >, < fda : 1 >$$

继续找出a的条件模式,继续和之前已经找到的比较,则余cefa1,cfda1

$$\begin{cases} b : < cfab : 1 >, < efc b : 1 >, < eac b : 1 > \\ c : < fdc : 1 > \\ a : < efa : 1 >, < fda : 1 > \end{cases}$$

综上,剩余条件模式为,因为同时满足min_sup和Item_sup所以X是长度为3的FT frequent

$$X = (bca)$$

Mining FT frequent itemsets of length equals to (δ+1)

Mining fault tolerant(FT) frequent itemsets using pattern growth:Design and construction

ABSTRACT

fault tolerant(FT) frequent itemset与普通的频繁项集模式挖掘更昂贵

Apriori系列算法会指数级的生成大量的候选项级,包含数据库没有的,并且会多次扫描数据库计算支持度

提出一种新颖的算法利用频繁模式增长方法挖掘频繁项级,采用自下而上的思想,将事务数据库中的投影到一个小数据库和事务数据库,通过事务数据库项级来挖掘频繁项级

FT-patternGrowth将事务数据库存储在一个高压压缩的小数据库和FT-tree中,直接在树中计数,无需多次扫描数据库,提高算法的处理速度

Introduction

普通的频繁项级挖掘方法在项目集部分缺失时会面临阈值设置问题,阈值过大挖掘过少,阈值过小会冗余

an itemset X with length greater than (δ+1) is a FT frequent itemsets if it has support of at least T number of FT-transactions

A transactions t is a FT-transaction of T under FT factor if it contains at least(|X|-δ) number of items of X

T is the support of X which must be greater or equals to the minimum itemset support

Each individual item i of X must appear in at least m number of FT-transactions of X

The m is the minimum item support under fault tolerant factor δ

Item	Support	Support
a	1	1
b	2	2
c	1	1
d	1	1
e	1	1
f	1	1

这里给出容错频繁挖掘的定义,假设min_supδ=3,Item_supδ=2,为一个不匹配,δ=1

FT-Apriori算法在阈值较小时处理时间长

项目集abcdef具有项目集支持为3,事务10,30,50包含五个项目中的4个,并且每个单独项目出现在三个事务中的两个中

主要基于Apriori的候选生成和测试,数据量较大

采用自下而上的搜索机制,在挖掘项集X之前先枚举了所有子集,这种机制限制了算法在合理的时间内挖掘完整的项级

在挖掘时会多次扫描数据库

消除了Apriori算法的候选生成测试和多次扫描数据库

第一次扫描数据库计算出所有长度为1的频繁项集,第二次扫描建立频繁模式树

Related work

Applications of mining FT frequent itemsets

从原数据库中挖掘频繁子图

使用频繁子图进行图像分类,图像的频繁近似子图被用于分类特征

确定和不确定的数据集的频繁模式挖掘算法

利用频繁项级挖掘生物数据库中的项级的算法,比例项集的相容性的错误数量与项集的长度成正比

基于Apriori算法,应用自下而上的完全搜索空间,应用向下封闭属性来剪枝却发现长度为k的子集不频繁则抛弃其所有超集

避免数据库重复扫描,只需扫描一次,为每个项目生成位向量,采用深度优先模式生成候选项集,使用位运算快速计算项级支持度,但是会生成大量的候选项集

VB-FT-Mine

pattern growth

construct more than ont FP-trees for each itemset to mine its supersets

在δ下,建立(δ+1)个FP-trees

在X和δ下,建立2^|X|个FP-trees

这两种方法都不能充分利用FP-tree的计数特性

挖掘比例项集的概念与传统的项集相似,但比例项集的容错因子与项集的长度成正比

$$I = i_1, \dots, i_m, X \in I \quad |X|$$

A transaction $T = (tid, t) \quad t = t_1, \dots, t_n$

$T \in TDB$

X is frequency $sup(X) \geq min_sup$

number of FT-transaction $\geq min_sup^\delta$

$m = item_sup^\delta$
Each individual item i of X must appear in at least m number of FT-transaction of X

Fault tolerant(FT) frequent itemset mining:Problem statement

容错频繁项集