



当前位置:解决方案网» 数据结构与算法 » FP-growth高效率频繁项集发现

FP-growth高效率频繁项集发现

本文收集于网络，只用于方便查找方案，感谢源作者，如果侵权请联系删除

FP-growth高效频繁项集发现

FP-growth

算法优缺点：

优点：一般快于**Apriori**

缺点：实现比较困难，在某些数据上性能下降

适用数据类型：标称型数据

算法思想：

FP-growth算法是用来解决频繁项集发现问题的，这个问题再前面我们可以通过**Apriori**算法来解决，但是虽然利用**Apriori**原理加快了速度，仍旧是效率比较低的。FP-growth算法则可以解决这个问题。

FP-growth算法使用了频繁模式树（Frequent Pattern Tree）的数据结构。FP-tree是一种特殊的前缀树，由频繁项头表和项前缀树构成。所谓前缀树，是一种存储候选项集的数据结构，树的分支用项名标识，树的节点存储后缀项，路径表示项集。

FP-growth算法生成频繁项集相对**Apriori**生成频繁项集的主要好处就是速度快，能快到几个数量级；另一个好处就是用FP树存储数据可以减少存储空间，因为关联挖掘的数据集往往是重复性很高的，这就能带来很高的压缩比。

算法可以分成一下几个部分：

构建FP树

首先我们要统计出所有的元素的频度，删除不满足最小支持度的（**Apriori**原理）
然后我们要根据频度对所有的项集排序(保证我们的树是最小的)
最后根据排序的项集构建FP树

从FP树挖掘频繁项集：

生成条件模式基
生成条件FP树

算法的执行过程这篇文章有个很好的示例程序

函数：

```
loadSimpDat()  
创建数据集  
createInitSet(dataSet)  
将数据集处理成字典的形式  
createTree(dataSet, minSup=1)  
创建FP树的主函数。首先生成单元元素的频繁项，然后对每个项集进行以频繁项的频度为基准的排序。  
updateTree(items, inTree, headerTable, count)  
根据每一个项集和对应的频数，更新FP树。并同时建立表头  
updateHeader(nodeToTest, targetNode)  
当指针已经初始化的时候，调用这个函数把新的点加到链表的最后面  
ascendTree(leafNode, prefixPath)  
向上遍历移植到根节点，将经过的节点都加到前缀路径中，得到整条每个频繁项的前缀路径  
findPrefixPath(basePat, treeNode)  
生成条件模式基  
mineTree(inTree, headerTable, minSup, preFix, freqItemList)  
递归调用生成条件FP树和频繁项集。创建条件FP树的过程可以重用前面createTree的代码  
1 #coding=utf-8
```

```

2 import time
3 class treeNode(object):
4     """docstring for treeNode"""
5     def __init__(self, nameValue, numOccur, parentNode):
6         super(treeNode, self).__init__()
7         self.name = nameValue
8         self.count = numOccur
9         self.nodeLink = None
10        self.parent = parentNode
11        self.children = {}
12    def inc(self, numOccur):
13        self.count += numOccur
14    def disp(self, ind=1):
15        print ' '*ind, self.name, ' ', self.count
16        for child in self.children.values():
17            child.disp(ind+1)
18    def loadSimpDat():
19        simpDat = [['r', 'z', 'h', 'j', 'p'],
20                    ['z', 'y', 'x', 'w', 'v', 'u', 't', 's'],
21                    ['z'],
22                    ['r', 'x', 'n', 'o', 's'],
23                    ['y', 'r', 'x', 'z', 'q', 't', 'p'],
24                    ['y', 'z', 'x', 'e', 'q', 's', 't', 'm']]
25        return simpDat
26    def createInitSet(dataSet):
27        retDict = {}
28        for trans in dataSet:
29            retDict[frozenset(trans)] = 1
30        return retDict
31    def createTree(dataSet, minSup=1):
32        headerTable = {}
33        #frequency of each item
34        for trans in dataSet:
35            for item in trans:
36                headerTable[item] = headerTable.get(item, 0) + dataSet[trans]#some trans may same
37        #remove items not meeting minSup
38        for k in headerTable.keys():
39            if headerTable[k] < minSup:
40                del(headerTable[k])
41        freqItemSet = set(headerTable.keys())
42        if len(freqItemSet) == 0:#no frequent item
43            return None, None
44        for k in headerTable:#add a point field
45            headerTable[k] = [headerTable[k], None]
46
47        retTree = treeNode('Null set', 1, None)
48        for tranSet, count in dataSet.items():
49            localD = {}
50            for item in tranSet:#把每一个项集的元素提取出来，并加上统计出来的频率
51                if item in freqItemSet:
52                    localD[item] = headerTable[item][0]
53            if len(localD) > 0:#排序，并更新树
54                orderdItem = [v[0] for v in sorted(localD.items(), key=lambda p:p[1], reverse=True)]
55                updateTree(orderdItem, retTree, headerTable, count)
56        return retTree, headerTable
57    def updateTree(items, inTree, headerTable, count):
58        #将新的节点加上来
59        if items[0] in inTree.children:
60            inTree.children[items[0]].inc(count)
61        else:
62            inTree.children[items[0]] = treeNode(items[0], count, inTree)
63            #更新指针
64            if headerTable[items[0]][1] == None:
65                headerTable[items[0]][1] = inTree.children[items[0]]
66            else:

```

```

67         updateHeader(headerTable[items[0]][1], inTree.children[items[0]])
68     if len(items) > 1:
69         updateTree(items[1:], inTree.children[items[0]], headerTable, count)
70 def updateHeader(nodeToTest, targetNode):
71     while nodeToTest.nodeLink != None:
72         nodeToTest = nodeToTest.nodeLink
73     nodeToTest.nodeLink = targetNode
74
75 def ascendTree(leafNode, prefixPath): #ascends from leaf node to root
76     if leafNode.parent != None:
77         prefixPath.append(leafNode.name)
78         ascendTree(leafNode.parent, prefixPath)
79
80 def findPrefixPath(basePat, treeNode): #treeNode comes from header table
81     condPats = {}
82     while treeNode != None:
83         prefixPath = []
84         ascendTree(treeNode, prefixPath)
85         if len(prefixPath) > 1:
86             condPats[frozenset(prefixPath[1:])] = treeNode.count
87         treeNode = treeNode.nodeLink
88     return condPats
89 def mineTree(inTree, headerTable, minSup, preFix, freqItemList):
90     bigL = [v[0] for v in sorted(headerTable.items(), key=lambda p: p[1])]#(sort header table)
91     #print bigL
92     for basePat in bigL: #start from bottom of header table
93         newFreqSet = preFix.copy()
94         newFreqSet.add(basePat)
95         print 'finalFrequent Item: ', newFreqSet #append to set
96         freqItemList.append(newFreqSet)
97         condPattBases = findPrefixPath(basePat, headerTable[basePat][1])
98         print 'condPattBases:', basePat, condPattBases
99         #2. construct cond FP-tree from cond. pattern base
100        myCondTree, myHead = createTree(condPattBases, minSup)
101        print 'head from conditional tree: ', myHead
102        if myHead != None: #3. mine cond. FP-tree
103            print 'conditional tree for: ', newFreqSet
104            myCondTree.disp(1)
105            mineTree(myCondTree, myHead, minSup, newFreqSet, freqItemList)
106 def main():
107     if True:
108         simpDat = loadSimpDat()
109         initSet = createInitSet(simpDat)
110         myFP, myHeadTable = createTree(initSet, 3)
111         myFP.disp()
112         freqItems = []
113         mineTree(myFP, myHeadTable, 3, set([]), freqItems)
114         print freqItems
115     if False:
116         t1 = time.clock()
117         parsedDat = [line.split() for line in open('kosarak.dat').readlines()]
118         initSet = createInitSet(parsedDat)
119         myFP, myHeadTable = createTree(initSet, 100000)
120         myfreq = []
121         mineTree(myFP, myHeadTable, 100000, set([]), myfreq)
122         t2 = time.clock()
123         print 'time=', t2-t1
124         print myfreq
125 if __name__ == '__main__':
126     main()
127

```

使用FP算法对一个近100万行的数据进行分析，耗时不过十来秒：

```
time= 12.5076159007
[set(['1']), set(['1', '6']), set(['3']), set(['11', '3']), set(['11', '3', '6']), set(['3', '6']), set(['11']), set(['11', '6']), set(['6'])]
[Finished in 13.2s]
```

www.data321.com 解决方案

而如果采用Apriori的频繁集发现算法我跑了四分多种没出结果然后就强制关掉了。。。

[Finished in 246.3s]

事实证明这个算法确实能够提高数量级的速度啊。

来自为知笔记(Wiz)

上一篇: 算法竞赛入门经典_第二章:循环结构程序设计_上机练习题_MyAnswer

下一篇: 《数据结构与算法分析:C语言描述_原书第二版》CH3表、栈跟队列_reading notes

猜你喜欢 **数据结构与算法** DBC:1

线性表的合并,该怎么处理 - 数据结构与算法

想要输入的字符串大于4096个字符怎么处理 - 数据结构与算法

数据结构与算法

数据结构与算法

数据结构与算法

acm有关问题 - 数据结构与算法

十个经典算法研究[新增GA、启发式搜索、SIFT、傅里叶变换等算法]解决方法 - 数据结

数据结构与算法

数据结构与算法

数据结构与算法

关于环形缓冲池的设计,该如何解决 - 数据结构与算法

二叉树创建有关问题,为什么输入一个字符后,一直输入'#'都结束不了(2) - 数据结构

0 条评论, 0 人参与。

★ 0



我有话说...

使用社交帐号登录

或以游客身份发布

昵称

最新评论

还没有评论

更多热评文章



serialVersionUID



手机页面是由HTML5写的



SSH整合开发后出现的org



架构最大的变化就是由Exchange

友言?



超快云服务器, 基础配置 **28.8** 元/月起

立刻抢购

本站声明: 本站所转载之内容, 无任何商业意图, 如涉及版权、著作权等问题, 请您告知: QQ:2450057300, 本站收到核实后会删除处理。
File:2015/12/6 4:40:44