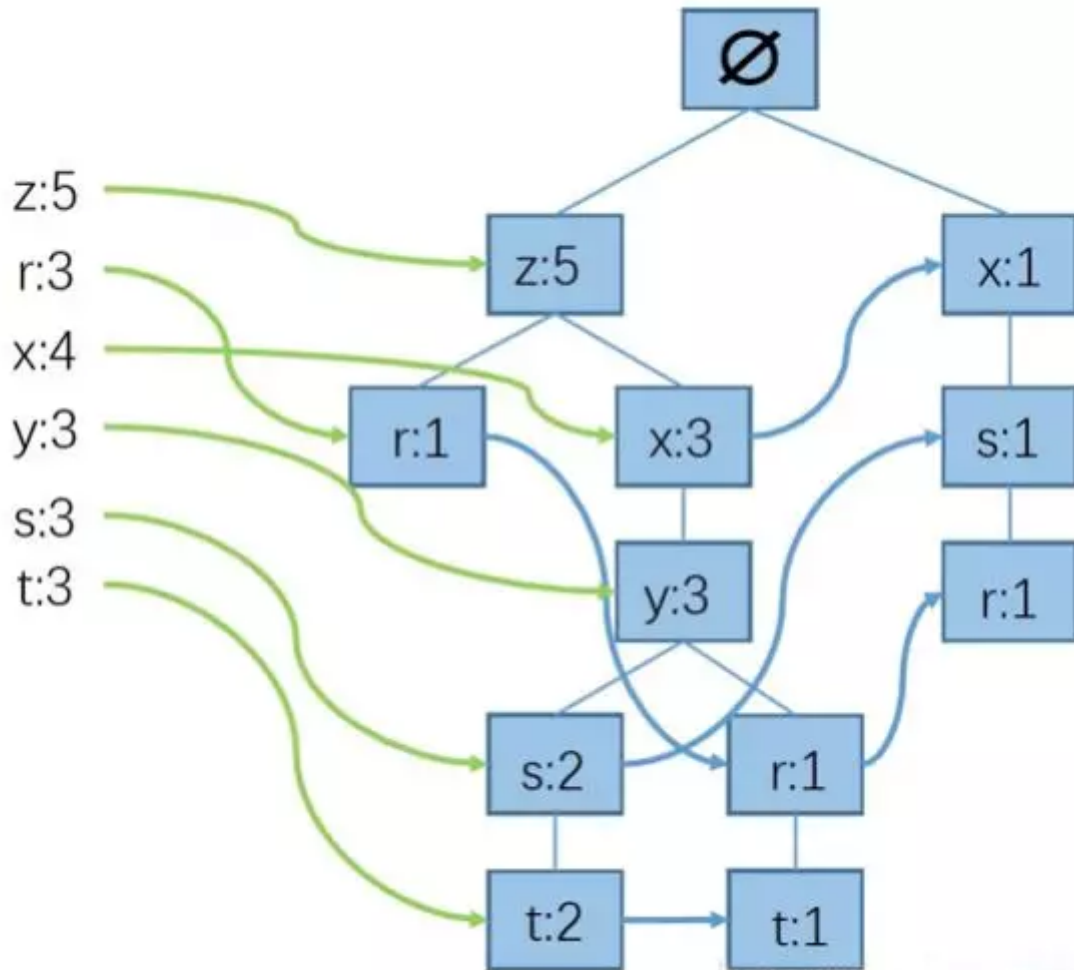# FP-growth算法的python实现

wsp　Python可视化编程机器学习OpenCV　2019-10-24

　　FP-growth算法是一种用于发现数据集中频繁模式的有效方法。Apriori算法在产生频繁模式完全集前需要对数据库进行多次扫描，同时产生大量的候选频繁集，这就使Apriori算法时间和空间复杂度较大。FP-growth算法由Apriori算法产生候选项集，然后扫描数据集来检查它们是否频繁。由于只对数据集扫描两次，因此它比Apriori算法速度要快，通常性能要好两个数量级以上。

　　在FP-growth算法中，数据集存储在一个称为FP(Frequent Pattern)树的结构中。FP树构建完成后，可以通过查找元素项的条件基以及构建条件FP树来发现频繁集。该过程不断以更多的元素为条件重复进行，知道FP树只包含一个元素为止。

　　下面仅以这个简单的数据集为例子--实际上，既使在多达百万条记录的大数据集上，FP-growth算法也能快速运行。

| instance id | elements |
| --- | --- |
| 0 | r, z, h, j, p |
| 1 | z, y, x, w, v, u, t, s |
| 2 | z |
| 3 | r, x, n, o, s |
| 4 | y, r, x, z, q, t, p |
| 5 | y, z, x, e, q, s, t, m |

python代码:

```
 1  '''
 2  FP-Growth FP means frequent pattern
 3  the FP-Growth algorithm needs:
 4  1. FP-tree (class treeNode)
 5  2. header table (use dict)
 6  This finds frequent itemsets similar to apriori but does not
 7  find association rules.
 8  @author: Peter
 9  '''
10  def loadSimpDat():
11      simpDat = [['r', 'z', 'h', 'j', 'p'],
12                 ['z', 'y', 'x', 'w', 'v', 'u', 't', 's'],
13                 ['z','p','x'],
14                 ['r', 'x', 'n', 'o', 's'],
15                 ['y', 'r', 'x', 'z', 'q', 't', 'p'],
16                 ['y', 'z', 'x', 'e', 'q', 's', 't', 'm']]
17      return simpDat
18
```

```
19  class treeNode:
20      def __init__(self, nameValue, numOccur, parentNode):
21          self.name = nameValue
22          self.count = numOccur
23          self.nodeLink = None
24          self.parent = parentNode       #needs to be updated
25          self.children = {}
26
27      def inc(self, numOccur):
28          self.count += numOccur
29
30      def disp(self, ind=1):
31          print (('  '*ind, self.name, ' ', self.count))
32          for child in self.children.values():
33              child.disp(ind+1)
34      #def __lt__(self, other):#定义 "<"用于sorted()
35          #return self.count < other.count
36
37  def createTree(dataSet, minSup=1): #create FP-tree from dataset but don't mi
38      headerTable = {}
39      #go over dataSet twice
40      for trans in dataSet:#first pass counts frequency of occurance
41          for item in trans:
42              headerTable[item] = headerTable.get(item, 0) + dataSet[trans]
43
44      for k in list(headerTable.keys()):  #remove items not meeting minSup
45          if headerTable[k] < minSup:
46              headerTable.pop(k)
47      freqItemSet = set(headerTable.keys())
48
49      #print 'freqItemSet: ',freqItemSet
50      if len(freqItemSet) == 0: return None, None  #if no items meet min suppo
51      for k in headerTable:
52          headerTable[k] = [headerTable[k], None] #reformat headerTable to use
53      #print 'headerTable: ',headerTable
54      retTree = treeNode('Null Set', 1, None) #create tree
55      for tranSet, count in dataSet.items():  #go through dataset 2nd time
56          localD = {}
57          for item in tranSet:  #put transaction items in order
58              if item in freqItemSet:
```

```python
            localD[item] = headerTable[item][0]
        if len(localD) > 0:
            orderedItems = [v[0] for v in sorted(localD.items(), key=lambda
            updateTree(orderedItems, retTree, headerTable, count)#populate
    return retTree, headerTable #return tree and header table

def updateTree(items, inTree, headerTable, count):
    if items[0] in inTree.children:#check if orderedItems[0] in retTree.chil
        inTree.children[items[0]].inc(count) #incrament count
    else:   #add items[0] to inTree.children
        inTree.children[items[0]] = treeNode(items[0], count, inTree)
        if headerTable[items[0]][1] == None: #update header table
            headerTable[items[0]][1] = inTree.children[items[0]]
        else:
            updateHeader(headerTable[items[0]][1], inTree.children[items[0]]
    if len(items) > 1:#call updateTree() with remaining ordered items
        updateTree(items[1::], inTree.children[items[0]], headerTable, count

def updateHeader(nodeToTest, targetNode):   #this version does not use recur
    while (nodeToTest.nodeLink != None):    #Do not use recursion to travers
        nodeToTest = nodeToTest.nodeLink
    nodeToTest.nodeLink = targetNode

def ascendTree(leafNode, prefixPath): #ascends from leaf node to root
    if leafNode.parent != None:
        prefixPath.append(leafNode.name)
        ascendTree(leafNode.parent, prefixPath)

def findPrefixPath(basePat, treeNode): #treeNode comes from header table
    condPats = {}
    while treeNode != None:
        prefixPath = []
        ascendTree(treeNode, prefixPath)
        if len(prefixPath) > 1:
            condPats[frozenset(prefixPath[1:])] = treeNode.count
        treeNode = treeNode.nodeLink
    return condPats

def mineTree(inTree, headerTable, minSup, preFix, freqItemList):
    bigL = [k for k,v in sorted(headerTable.items(), key=lambda p: p[1][0])]
```

```python
 99         for basePat in bigL:  #start from bottom of header table
100             newFreqSet = preFix.copy()
101             newFreqSet.add(basePat)
102             #print 'finalFrequent Item: ',newFreqSet    #append to set
103             freqItemList.append(newFreqSet)
104             condPattBases = findPrefixPath(basePat, headerTable[basePat][1])
105             #print 'condPattBases :',basePat, condPattBases
106             #2. construct cond FP-tree from cond. pattern base
107             myCondTree, myHead = createTree(condPattBases, minSup)
108             #print 'head from conditional tree: ', myHead
109             if myHead != None: #3. mine cond. FP-tree
110                 #print 'conditional tree for: ',newFreqSet
111                 #myCondTree.disp(1)
112                 mineTree(myCondTree, myHead, minSup, newFreqSet, freqItemList)#,
113
114 def createInitSet(dataSet):
115     retDict = {}
116     for trans in dataSet:
117         retDict[frozenset(trans)] = 1
118     return retDict
119
120
121 minSup = 4
122 simpDat = loadSimpDat()
123 initSet = createInitSet(simpDat)
124 myFPtree, myHeaderTab = createTree(initSet, minSup)
125 myFreqList = []
126 if myFPtree is not None:
127     myFPtree.disp()
128     mineTree(myFPtree, myHeaderTab, minSup, set([]), myFreqList)
129 print("支持度为%d时，频繁项数为%d:"%(minSup, len(myFreqList)))
130 print("频繁项集为:")
131 for item in myFreqList:
132     print(item)
```

```
(' ', 'Null Set', '', 1)
('   ', 'z', '', 2)
('    ', 'x', '', 1)
('    ', 'x', '', 4)
('     ', 'z', '', 3)
支持度为4时，频繁项数为2:
频繁项集为:
{'z'}
{'x'}
```

喜欢此内容的人还喜欢

## Android 入门程序 Kotlin版（1）

Python可视化编程机器学习OpenCV

---

## 回老家要隔离、做核酸吗？用它查

腾讯

---

## 成熟的人，要学会从"负重前行"，到"举重若轻"

刘润