# 机器学习学习报告

周珊琳

上海电力大学
上海, 2019-06-28, 中国

## 第一节　算法概述

算法名称：FP-tree

算法原理：FP-growth算法不同于Apriori算法生成候选项集再检验是否频繁的"产生-测试"方法，而是使用一种称为频繁模式树（FP-tree）的紧凑数据结构组织数据，并直接从该结构中提取频繁项集。FP-树是一种输入数据的压缩表示，通过逐个读入事务，并把每个事务映射到FP-树中的一条路径来构造。由于不同的事务可能会有若干个相同的项，因此它们的路径可能部分重叠。路径相互重叠越多，使用FP-tree结构获得的压缩效果越好。如果FP-tree足够小，能够存放在内存中，就可以直接从这个内存中的结构提取频繁项集，而不必重复地扫描放在硬盘上的数据，从而提高处理效率。

算法中所用的到数学公式

支持度公式：

$$Support(A \rightarrow B) = support_count(A \bigcup B)/N \tag{1.1}$$

置信度公式：

$$confidence(A \rightarrow B) = support_count(A \bigcup B)/support_count(A) \tag{1.2}$$

# 第二节　算法设计

## 2.1　算法流程

---
**Algorithm 1** FP-growth 算法
---
1: 收集数据
2: 准备数据：由于存储的是集合，需要离散数据，即连续型数据需要离散化
3: 分析数据
4: 训练算法：构建一个FP树，对其进行挖掘（抽取频繁项集)
5: 使用算法

---

---
**Algorithm 2** FP-growth 抽取频繁项集的步骤:
---
1: 从FP中获取条件模式基
2: 利用条件模式基，构建一个条件FP树
3: 迭代重复前两个步骤，直到树包含一个元素为止

---

## 2.2　核心代码

源代码 1:

```
1
2  class treeNode:
3      def __init__(self, nameValue, numOccur, parentNode):
4          self.name =nameValue
5          self.count =numOccur
6          self.nodeLink =None
7          self.parent =parentNode      #needs to be updated
8          self.children ={}
9
10     def inc(self, numOccur):
11         self.count +=numOccur
12
13     def disp(self, ind=1):
14         print( ' '*ind, self.name, ' ', self.count)
15         for child in self.children.values():
16             child.disp(ind+1)
17
18 def createTree(dataSet, minSup=1): #create FP-tree from dataset but don't mine
19     headerTable ={}
20     #go over dataSet twice
21     for trans in dataSet:#first pass counts frequency of occurance
22         for item in trans:
23             headerTable[item] =headerTable.get(item, 0) +dataSet[trans]
24     for k in headerTable.keys(): #remove items not meeting minSup
25         if headerTable[k] <minSup:
26             del(headerTable[k])
27     freqItemSet =set(headerTable.keys())
```

```python
28        #print 'freqItemSet: ',freqItemSet
29        if len(freqItemSet) ==0: return None, None #if no items meet min support -->get out
30        for k in headerTable:
31            headerTable[k] =[headerTable[k], None] #reformat headerTable to use Node link
32        #print 'headerTable: ',headerTable
33        retTree =treeNode('Null Set', 1, None) #create tree
34        for tranSet, count in dataSet.items(): #go through dataset 2nd time
35            localD ={}
36            for item in tranSet: #put transaction items in order
37                if item in freqItemSet:
38                    localD[item] =headerTable[item][0]
39            if len(localD) >0:
40                orderedItems =[v[0] for v in sorted(localD.items(), key=lambda p: p[1], reverse=True)]
41                updateTree(orderedItems, retTree, headerTable, count)#populate tree with ordered freq
                                                    itemset
42        return retTree, headerTable #return tree and header table
43
44  def updateTree(items, inTree, headerTable, count):
45      if items[0] in inTree.children:#check if orderedItems[0] in retTree.children
46          inTree.children[items[0]].inc(count) #incrament count
47      else:  #add items[0] to inTree.children
48          inTree.children[items[0]] =treeNode(items[0], count, inTree)
49          if headerTable[items[0]][1] ==None: #update header table
50              headerTable[items[0]][1] =inTree.children[items[0]]
51          else:
52              updateHeader(headerTable[items[0]][1], inTree.children[items[0]])
53      if len(items) >1:#call updateTree() with remaining ordered items
54          updateTree(items[1::], inTree.children[items[0]], headerTable, count)
55
56  def updateHeader(nodeToTest, targetNode): #this version does not use recursion
57      while (nodeToTest.nodeLink !=None): #Do not use recursion to traverse a linked list!
58          nodeToTest =nodeToTest.nodeLink
59      nodeToTest.nodeLink =targetNode
60
61  def ascendTree(leafNode, prefixPath): #ascends from leaf node to root
62      if leafNode.parent !=None:
63          prefixPath.append(leafNode.name)
64          ascendTree(leafNode.parent, prefixPath)
65
66  def findPrefixPath(basePat, treeNode): #treeNode comes from header table
67      condPats ={}
68      while treeNode !=None:
69          prefixPath =[]
70          ascendTree(treeNode, prefixPath)
71          if len(prefixPath) >1:
72              condPats[frozenset(prefixPath[1:])] =treeNode.count
73          treeNode =treeNode.nodeLink
74      return condPats
75
76
77  def mineTree(inTree, headerTable, minSup, preFix, freqItemList):
78      bigL =[v[0] for v in sorted(headerTable.items(), key=lambda p: p[1])]#(sort header table)
79      for basePat in bigL: #start from bottom of header table
80          newFreqSet =preFix.copy()
```

```
81          newFreqSet.add(basePat)
82          #print 'finalFrequent Item: ',newFreqSet #append to set
83          freqItemList.append(newFreqSet)
84          condPattBases =findPrefixPath(basePat, headerTable[basePat][1])
85          #print 'condPattBases :',basePat, condPattBases
86          #2. construct cond FP-tree from cond. pattern base
87          myCondTree, myHead =createTree(condPattBases, minSup)
88          #print 'head from conditional tree: ', myHead
89          if myHead !=None: #3. mine cond. FP-tree
90              #print 'conditional tree for: ',newFreqSet
91              #myCondTree.disp(1)
92              mineTree(myCondTree, myHead, minSup, newFreqSet, freqItemList)
93
94  def loadSimpDat():
95      simpDat =[['r', 'z', 'h', 'j', 'p'],
96                ['z', 'y', 'x', 'w', 'v', 'u', 't', 's'],
97                ['z'],
98                ['r', 'x', 'n', 'o', 's'],
99                ['y', 'r', 'x', 'z', 'q', 't', 'p'],
100               ['y', 'z', 'x', 'e', 'q', 's', 't', 'm']]
101     return simpDat
102
103 def createInitSet(dataSet):
104     retDict ={}
105     for trans in dataSet:
106         retDict[frozenset(trans)] =1
107     return retDict
108
109 import twitter
110 from time import sleep
111 import re
112
113 def textParse(bigString):
114     urlsRemoved =re.sub('(http:[/][/]|www.)([a-z]|[A-Z]|[0-9]|[/.]|[~])*', '', bigString)
115     listOfTokens =re.split(r'\W*', urlsRemoved)
116     return [tok.lower() for tok in listOfTokens if len(tok) >2]
117
118 def getLotsOfTweets(searchStr):
119     CONSUMER_KEY =''
120     CONSUMER_SECRET =''
121     ACCESS_TOKEN_KEY =''
122     ACCESS_TOKEN_SECRET =''
123     api =twitter.Api(consumer_key=CONSUMER_KEY, consumer_secret=CONSUMER_SECRET,
124                      access_token_key=ACCESS_TOKEN_KEY,
125                      access_token_secret=ACCESS_TOKEN_SECRET)
126     #you can get 1500 results 15 pages * 100 per page
127     resultsPages =[]
128     for i in range(1,15):
129         print( "fetching page %d" % i)
130         searchResults =api.GetSearch(searchStr, per_page=100, page=i)
131         resultsPages.append(searchResults)
132         sleep(6)
133     return resultsPages
134
```

```python
135  def mineTweets(tweetArr, minSup=5):
136      parsedList =[]
137      for i in range(14):
138          for j in range(100):
139              parsedList.append(textParse(tweetArr[i][j].text))
140      initSet =createInitSet(parsedList)
141      myFPtree, myHeaderTab =createTree(initSet, minSup)
142      myFreqList =[]
143      mineTree(myFPtree, myHeaderTab, minSup, set([]), myFreqList)
144      return myFreqList
145
146  from numpy import *
147  import numpy as np
148  from sklearn.datasets import load_iris,load_wine
149  import pandas as pd
150
151  #iris
152  #data=load_iris()
153  #dataSet = data.data
154  #labels = data.feature_names
155  #target = data.target
156  #r=[]
157  #for i in target:
158  #    if i==0:
159  #        i='Iris-setosa'
160  #    elif i==1:
161  #        i='Iris-versicolor'
162  #    else:
163  #        i='Iris-virginica'
164  #    r.append(i)
165  #simpDat = []
166  #for i in range(4):
167  #    tmp=pd.cut(dataSet[:,1],3,labels=['a'+str(i),'b'+str(i),'c'+str(i)])
168  #    simpDat.append(tmp)
169  #simpDat.append(r)
170  #simpDat = np.array(list(zip(*simpDat)))
171
172  #wine
173  #data=load_wine()
174  #dataSet = data.data
175  #r = data.target
176  #simpDat = []
177  #for i in range(13):
178  #    tmp=pd.cut(dataSet[:,1],3,labels=['a'+str(i),'b'+str(i),'c'+str(i)])
179  #    simpDat.append(tmp)
180  #simpDat.append(r)
181  #simpDat = np.array(list(zip(*simpDat)))
182
183  def loadDataSet(fileName):
184      dataMat =[]
185      fr = open(fileName)
186      for line in fr.readlines():
187          lineArr =line.strip().split(',')
188          # print(lineArr)
```

```
189          if lineArr[0] =='vhigh':
190              lineArr[0] ='a1'
191          if lineArr[0] =='high':
192              lineArr[0] ='a2'
193          if lineArr[0] =='med':
194              lineArr[0] ='a3'
195          if lineArr[0] =='low':
196              lineArr[0] ='a4'
197
198          if lineArr[1] =='vhigh':
199              lineArr[1] ='b1'
200          if lineArr[1] =='high':
201              lineArr[1] ='b2'
202          if lineArr[1] =='med':
203              lineArr[1] ='b3'
204          if lineArr[1] =='low':
205              lineArr[1] ='b4'
206
207          if lineArr[2]=='2':
208              lineArr[2]='c1'
209          if lineArr[2]=='3':
210              lineArr[2]='c2'
211          if lineArr[2]=='4':
212              lineArr[2]='c3'
213          if lineArr[2]=='5more':
214              lineArr[2]='c4'
215
216          if lineArr[3]=='2':
217              lineArr[3]='d1'
218          if lineArr[3]=='4':
219              lineArr[3]='d2'
220          if lineArr[3]=='more':
221              lineArr[3]='d3'
222
223          if lineArr[4]=='small':
224              lineArr[4]='e1'
225          if lineArr[4]=='med':
226              lineArr[4]='e2'
227          if lineArr[4]=='big':
228              lineArr[4]='e3'
229
230          if lineArr[5]=='low':
231              lineArr[5]='f1'
232          if lineArr[5]=='med':
233              lineArr[5]='f2'
234          if lineArr[5]=='high':
235              lineArr[5]='f3'
236
237          dataMat.append(lineArr)
238
239      return dataMat
240  #car
241  simpDat =loadDataSet('car.data')
242
```

```
243
244  #initSet = createInitSet(simpDat)
245  #print(initSet)
246  #myFPtree, myHeaderTab = createTree(initSet, minSup)
247  #myFPtree.disp()
248  #myFreqList = []
249  #mineTree(myFPtree, myHeaderTab, minSup, set([]), myFreqList)
250  #print(myFreqList)
251
252  minSup=300
253  #simpDat = loadSimpDat()
254  #print(simpDat)
255  initSet =createInitSet(simpDat)
256  #print(initSet)
257  myFPtree, myHeaderTab =createTree(initSet, minSup)
258  myFPtree.disp()
259  myFreqList =[]
260  mineTree(myFPtree, myHeaderTab, minSup, set([]), myFreqList)
261  print(myFreqList)
262
263  #parsedDat = [line.split() for line in open('kosarak.dat').readlines()]
264  #initSet = createInitSet(parsedDat)
265  #myFPtree, myHeaderTab = createTree(initSet, 100000)
266  #myFPtree.disp()
267  #myFreqList = []
268  #mineTree(myFPtree, myHeaderTab, 100000, set([]), myFreqList)
269  #print(myFreqList)
```

# 第三节　选用数据

iris行数：150 列数：5
列属性及取值:
1)萼片长度cm，数值型
2)萼片宽度cm，数值型
3)花瓣长度cm，数值型
4)花瓣宽度cm数值型
类别：
Iris Setosa
Iris Versicolour
Iris Virginica
car，行数：1728，列数：6
列属性及取值：
1)buying: vhigh, high, med, low.
2)maint: vhigh, high, med, low.
3)doors: 2, 3, 4, 5more.

4)persons: 2, 4, more.

5)lugboot: small, med, big.

6)safety: low, med, high.

类别：

unacc, acc, good, vgood

wine，行数：178，列数：13

属性：

1) Alcohol

2) Malic acid

3) Ash

4) Alcalinity of ash

5) Magnesium

6) Total phenols

7) Flavanoids

8) Nonflavanoid phenols

9) Proanthocyanins

10) Color intensity

11) Hue

12) OD280/OD315 of diluted wines

13) Proline

类别：

Alcohol 1，2，3

news新闻网站点击流：

列属性：

新闻报道编号

# 第四节　实验结果截图

# 第五节　实验分析和比较

对于iris，wine来说使用该算法的意义不大，只能找出哪类超过了3+；相对于前两个数据集car本身是离散型的数据，所以结果还行，找出来了哪类车超过了500+。因为car原来的不同属性列直接有相同的取值，因此需要事先进行处理，而iris和wine因为是连续型数据，因此需要提前离散化处理。后选择了一个比较适合这个算法的新闻网站点击流，找出有哪些新闻报道曾被10万+人浏览过。

实验分析和比较



图 1: iris数据展示
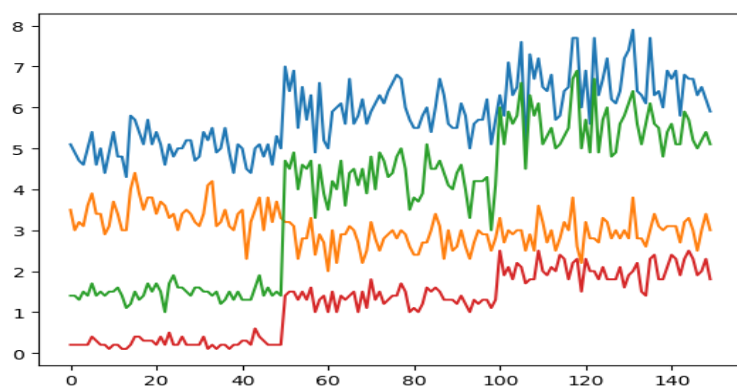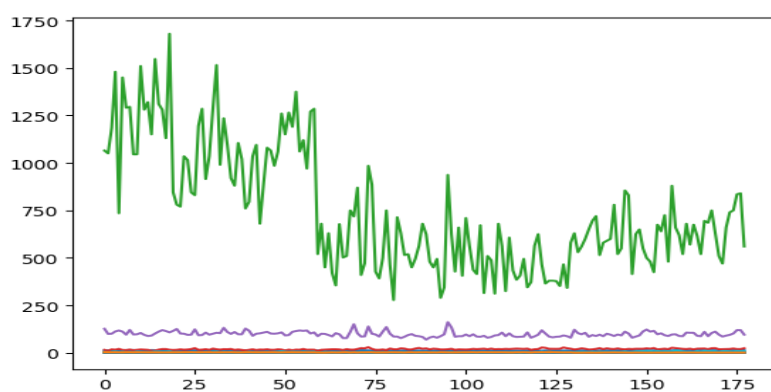


图 2: wine数据展示

```
Null Set   1
   Iris-setosa    3
      b1   1
        b0   1
          b3   1
            b2   1
      a0   1
        a3   1
          a1   1
            a2   1
    b1   2
      b0   1
        b3   1
          b2   1
      Iris-virginica    1
        b0   1
          b3   1
            b2   1
    a0   1
      a3   1
        a1   1
          a2   1
    Iris-virginica    2
      a0   1
        a3   1
          a1   1
            a2   1
[{'Iris-setosa'}, {'b1'}, {'b0'}, {'b0', 'b1'}, {'b3'}, {'b3', 'b0'}, {'b3', 'b1'}, {'b3', 'b0', 'b1'},
{'b2'}, {'b2', 'b3'}, {'b2', 'b0'}, {'b2', 'b0', 'b3'}, {'b2', 'b1'}, {'b2', 'b0', 'b1'}, {'b2', 'b3',
'b1'}, {'b2', 'b3', 'b0', 'b1'}, {'a0'}, {'a3'}, {'a3', 'a0'}, {'a1'}, {'a1', 'a3'}, {'a1', 'a0'}, {'a1',
'a3', 'a0'}, {'a2'}, {'a1', 'a2'}, {'a2', 'a3'}, {'a1', 'a2', 'a3'}, {'a2', 'a0'}, {'a1', 'a2', 'a0'},
{'a2', 'a3', 'a0'}, {'a1', 'a2', 'a3', 'a0'}, {'Iris-virginica'}]
```

图 3: iris运行结果

```
{'a10', 'a12', 'a3', 'a5', 'a6', 'a7', 'a9'},
{'a10', 'a12', 'a3', 'a4', 'a5', 'a6', 'a7', 'a9'},
{'a11', 'a12', 'a3'},
{'a11', 'a12', 'a3', 'a7'},
{'a11', 'a12', 'a3', 'a4'},
{'a11', 'a12', 'a3', 'a4', 'a7'},
{'a11', 'a12', 'a3', 'a5'},
{'a11', 'a12', 'a3', 'a4', 'a5'},
{'a11', 'a12', 'a3', 'a5', 'a7'},
{'a11', 'a12', 'a3', 'a4', 'a5', 'a7'},
{'a10', 'a11', 'a12', 'a3'},
{'a10', 'a11', 'a12', 'a3', 'a4'},
{'a10', 'a11', 'a12', 'a3', 'a5'},
{'a10', 'a11', 'a12', 'a3', 'a4', 'a5'},
{'a10', 'a11', 'a12', 'a3', 'a7'},
{'a10', 'a11', 'a12', 'a3', 'a4', 'a7'},
{'a10', 'a11', 'a12', 'a3', 'a5', 'a7'},
{'a10', 'a11', 'a12', 'a3', 'a4', 'a5', 'a7'},
{'a11', 'a12', 'a3', 'a9'},
{'a11', 'a12', 'a3', 'a4', 'a9'},
{'a11', 'a12', 'a3', 'a5', 'a9'},
{'a11', 'a12', 'a3', 'a4', 'a5', 'a9'},
{'a11', 'a12', 'a3', 'a7', 'a9'},
{'a11', 'a12', 'a3', 'a4', 'a7', 'a9'},
{'a11', 'a12', 'a3', 'a5', 'a7', 'a9'},
{'a11', 'a12', 'a3', 'a4', 'a5', 'a7', 'a9'},
{'a10', 'a11', 'a12', 'a3', 'a9'},
{'a10', 'a11', 'a12', 'a3', 'a4', 'a9'},
{'a10', 'a11', 'a12', 'a3', 'a7', 'a9'},
{'a10', 'a11', 'a12', 'a3', 'a4', 'a7', 'a9'},
{'a10', 'a11', 'a12', 'a3', 'a5', 'a9'},
{'a10', 'a11', 'a12', 'a3', 'a4', 'a5', 'a9'},
{'a10', 'a11', 'a12', 'a3', 'a5', 'a7', 'a9'},
```

图 4: wine运行结果

```
        d3   48

[{'d2'}, {'d3'}, {'f3'}, {'e3'}, {'e1'}, {'f1'}, {'f1', 'unacc'}, {'f2'}, {'e2'}, {'d1'}, {'unacc',
'd1'}, {'unacc'}]


    [''']
```

图 5: car运行结果

```
Null Set    1
    3   76514
      1   12917
    1   16829
    6   412762
      11   261773
        3   117401
          1   34141
          1   43366
        3   68888
          1   13436
          1   16461
      11   21190
        3   9718
          1   1565
          1   1882
[{'1'}, {'6', '1'}, {'3'}, {'11', '3'}, {'6', '11', '3'}, {'6', '3'}, {'11'}, {'6', '11'}, {'6'}]
```

图 6: news运行结果

10

## 第六节　遇到的问题及解决方法，实践心得

　　一开始对如何自底向上回溯找频繁项集不知道如何实现，后来参考了《机器学习实战》。还有就是对数据的处理，尤其是car，后来才发现car存在不同属性列直接有相同的取值这个问题，所以需要进行处理。每个数据都有其适合的算法，有时候硬套算法，可能意义不大。