

实验报告

141220132 银琦 141220132@smail.nju.edu.cn

一. 实验要求

By varying different level of support and confidence, compare Apriori, FPGrowth and a dummy the baseline method that conducting exhaustive search for frequent itemsets, in terms of the number of generated frequent itemsets, the storage used for mining the rules, the computational cost (in second)

Try to discover some interesting association rules using Apriori or FPGrowth, make discussion on the insights the rules disclose.

二. 问题及数据集描述

给定数据集，挖掘出数据集中的频繁项集，通过调整支持度与置信度便可得出关联规则。

给定了两个数据集，第一个是杂货店一个月的交易记录，每行代表一条交易；第二个数据集是 8 个用户在 9 台计算机的 Unix 系统上两年内使用的指令，每组指令都由**SOF**和**EOF**分隔开。

三. 方法描述

本次实验共比较了三种方法对数据的处理：Apriori 算法、FPGrowth 算法和蛮力算法。

1. Apriori 算法

算法使用频繁项集性质的先验知识，通过连接产生候选选项及其支持度，然后通过剪枝生成频繁项集。对于该种算法及给定的数据集，在 Weka 中跑不出结果，因为消耗大量的内存和时间，所以在 matlab 中手动实现。

2. FPGrowth 算法

针对 Apriori 算法固有的多次扫描事务数据集的缺陷，首先将代表频繁项集的数据库压缩到一颗频繁模式树，该树仍保留项集的关联信息，然后把这种压缩后的数据库划分成一组条件数据库，每个数据库关联一个频繁项或“模式段”，并分别挖掘每个条件数据库，对于每个“模式片段”只需要考察与它相关联数据集。对于该种算法及给定数据集，在 Weka 中可得出第一个数据集的关联规则，未手动实现。

3. 蛮力算法

未实现。

四. 方法的比较

1. Apriori 算法

用 Java 语言对数据集做了处理，将原数据集文件中的数据提取出，存入 txt 文件

中，数据集一格式如下：

```
1 citrus fruit,semi-finished bread,margarine,ready soups
2 tropical fruit,yogurt,coffee
3 whole milk
4 pip fruit,yogurt,cream cheese ,meat spreads
5 other vegetables,whole milk,condensed milk,long life bakery product
6 whole milk,butter,yogurt,rice,abrasive cleaner
7 rolls/buns
8 other vegetables,UHT-milk,rolls/buns,bottled beer,liquor (appetizer)
9 pot plants
10 whole milk,cereals
11 tropical fruit,other vegetables,white bread,bottled water,chocolate
12 citrus fruit,tropical fruit,whole milk,butter,curd,yogurt,flour,bottled water,dishes
13 beef
14 frankfurter,rolls/buns,soda
15 chicken,tropical fruit
16 butter,sugar,fruit/vegetable juice,newspapers
17 fruit/vegetable juice
18 packaged fruit/vegetables
19 chocolate
20 specialty bar
21 other vegetables
22 butter milk,pastry
23 whole milk
24 tropical fruit,cream cheese ,processed cheese,detergent,newspapers
25 tropical fruit,root vegetables,other vegetables,frozen dessert,rolls/buns,flour,sweet spreads,salty snack,waffles,candy,bathroom cleaner
26 bottled water,canned beer
27 yogurt
28 sausage,rolls/buns,soda,chocolate
29 other vegetables
30 brown bread,soda,fruit/vegetable juice,canned beer,newspapers,shopping bags
31 yogurt,beverages,bottled water,specialty bar
32 hamburger meat,other vegetables,rolls/buns,spices,bottled water,hygiene articles,napkins
33 root vegetables,other vegetables,whole milk,beverages,sugar
34 pork,berries,other vegetables,whole milk,whipped/sour cream,artif. sweetener,soda,abrasive cleaner
35 beef,grapes,detergent
36 pastry,soda
37 fruit/vegetable juice
38 canned beer
39 root vegetables,other vegetables,whole milk,dessert
40 citrus fruit,zwieback,newspapers
41 sausage,rolls/buns,soda,canned beer,specialty bar,shopping bags
42 tropical fruit,root vegetables,whole milk,yogurt,domestic eggs,brown bread,pastry,sugar,cereals,coffee,soda,waffles,candy
43 berries,yogurt
44 canned beer
```

数据集二格式如下：

```
1 whoami, pwd, ls, dir, vi, source, source, exit
2 whereis, mkdir, vi, vi, ls, source
3 elm, log
4 elm, man, fg, finger, |, more, finger, fg
5 elm, logout
6 elm, cd, ls, vi, fg, exit
7 elm, logout
8 elm, exit
9 elm, log
10 cd, cd, ls, vi, vi, exit
11 elm, exit
12 lwinfo, lwinfo, lwinfo, lwinfo, lwinfo, lwinfo, exit
13 elm
14 elm, exit
15 elm, cd, ls, vi, fg, vi, fg, cd, elm, finger, finger, |, more, fg, finger, finger, |, more, fg, finger, |, more, fg,
16 ls, more, exit
17 elm, finger, fg, finger, fg, finger, fg, finger, finger, finger, |, more, fg, finger, |, more, date, fg, elm, ls, fg, ls, emacs, ls, ls, emacs, ls, cat, rm
18 cd, ls, vi, exit
19 ls, cd, ls, mkdir, cd, ls, vi, vi, exit
20 exit
21 ls, uptime, exit
22 ruptime, exit
23 elm, vi, more, fg, finger, fg, ls, cd, ls, cd, ls, cd, vi, ls, more, ls, ls, -a, cd, ls, ls, -R, |, grep, finger, talk, fg, exit
24 whoami, ls, uptime, netscape
25 who, exit
26 elm, uptime, elm, r
27 setenv, ls, more, rm, ls, cd, ls, cd, ls, maker, ruptime
28 elm, finger, ruptime, ruptime, |, more, home, fg, elm
29 elm, finger, finger, fg, finger, finger, finger
30 elm, logout
31 ls, cd, ls, more, elm, exit
32 elm, exit
33 elm, ls, cd, ls, vi, ls, wc, ls, wc, mv, ls, rm, exit, cd, fg, exit
34 ls, ls, cd, ls, vi, ls, mv, mv, ll, cd, whereis, <GENSYM:0>/vi, whereis, vi, exit
35 cd, cd, ls, more, ls, vi, <GENSYM:0>/vi, ls, rm, exit
36 elm, ls, ll, fg, elm, exit
37 cd, ls, vi, cls, man, cls, ls, vi, ls
38 write, exit
39 setenv, cd, ls, maker, ls, cd, ls, cd, ls, cd, ls, ghostview, ls, cd, ls, cd, ls, man, ls, lp, -dmsl21, -n15, lpq, lpq, -dmsl21, lpq, -Pmsl21, lpq, -dmsl21
40 elm, log
41 cancel, -u, cancel, -u, ls, cd, lp, -dmsl21, cancel, -u, lp, -d, cancel, -u, lwifo, lwinfo, ls, lwinfo, lwinfo, back, back, cd, ls, vi, ls, <GENSYM:0>/vi,
42 setenv, cd, ls, vi, ls, vi, ls, vi, ls, lwinfo, lwinfo, lwinfo, lwinfo, ls, cancel, -u, lp, -dmsl22, lp, -dmsl22, lpq, -Pmsl22, lpq, -Pmsl23, cancel, -u, l
43 elm, lwinfo, lwinfo, lwinfo, lwinfo, lwinfo, lpq, -Pmsl23, cancel, -u, cancel, -u, uclln, lpq, -Pmsl22, lpq, -Pmsl23, cancel, -u, cancel, -u, man, ls, lpq, -P
44 setenv, ls, cd, maker, ls, ls, ls, ls, vi
```

在 matlab 中读入文件，将其转化成 01 矩阵，用 Apriori 算法得出关联规则，一开始设置支持度为 0.001，置信度为 0.9，但是运行时间过长，没有得出结果，之后将支持度改为 0.01，置信度改为 0.5，得出了结果。

2. FPGrowth 算法

用 Java 语言对数据集做了处理，将原数据集文件中的数据提取出，存入 arff 文件中，数据集一格式如下：

[illegible]

数据集二格式如下:

[illegible]

用 Weka 打开 arff 文件, 选择 FPGrowth 算法进行挖掘, 对数据集一分别测试了支持度为 0.001, 置信度为 0.9 和支持度为 0.01, 置信度为 0.5 的情况并得出结果, 对数据集二仅测试了支持度为 0.01, 置信度为 0.5 的情况并得出结果。

五. 结果分析

1. Apriori 算法

```
>> apriori
关联规则算法完成, 规则数为: 14
存储规则到文件 'rules.txt' 完成
Apriori算法挖掘菜品订单关联规则完成!
用时:
```

```
ans =
```

```
35.9590
```

```
>> apriori
关联规则算法完成, 规则数为: 100
存储规则到文件 'rules2.txt' 完成
Apriori算法挖掘菜品订单关联规则完成!
用时:
```

```
ans =
```

```
88.8780
```

由上图可以得出, 当支持度为 0.01, 置信度为 0.5 时, 该算法对于数据集一可以得出 14 条关联规则, 耗时大约 36 秒, 对数据集二可得出 100 条规则, 耗时大约 89 秒。

2. FPGrowth 算法

Associator

Choose **FPGrowth** -P 2 -I -1 -N 20 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.001

Start Stop

Associator output

Scheme: weka.associations.FPGrowth -P 2 -I -1 -N 20 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.001

Relation: Groceries

Instances: 9835

Attributes: 169

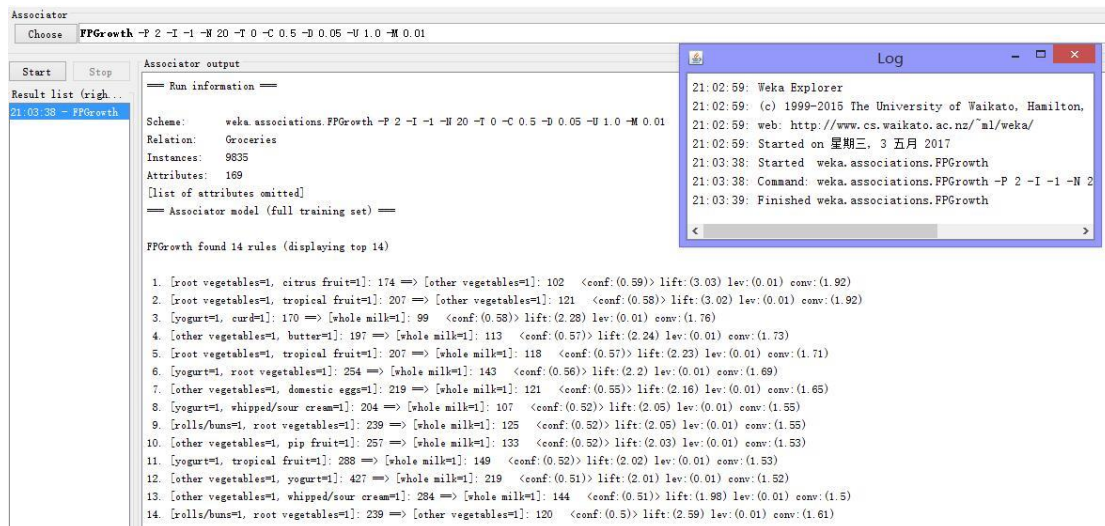
[List of attributes omitted]

== Associator model (full training set) ==

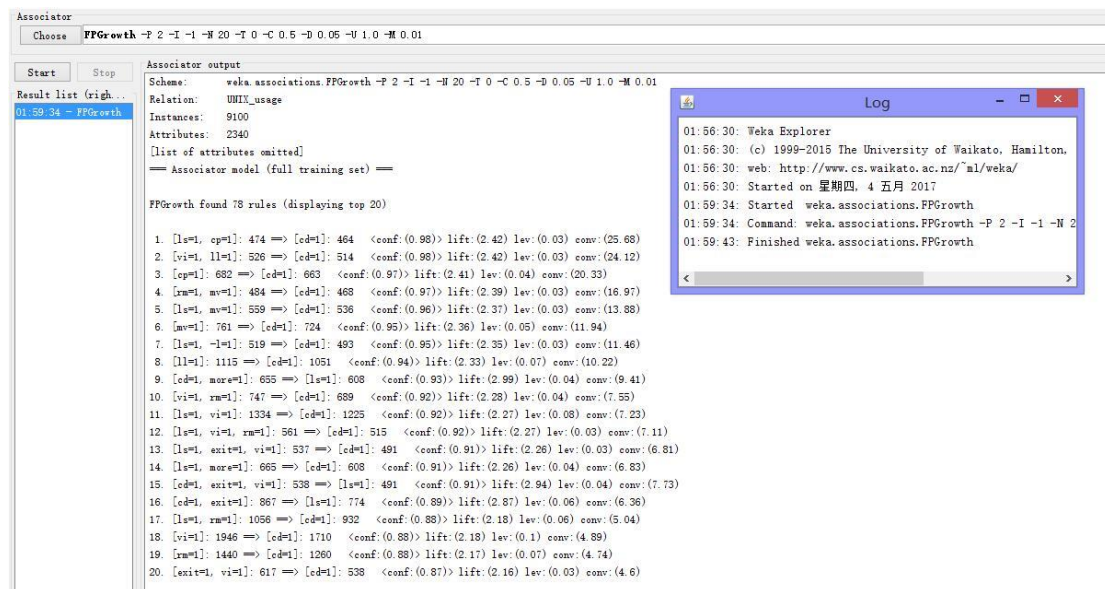
FPGrowth found 128 rules (displaying top 20)

1. [sugar=1, rice=1]: 12 ==> [whole milk=1]: 12 <conf:(1)> lift:(3.91) lev:(0) conv:(8.93)
2. [hygiene articles=1, canned fish=1]: 11 ==> [whole milk=1]: 11 <conf:(1)> lift:(3.91) lev:(0) conv:(8.19)
3. [root vegetables=1, pip fruit=1, hygiene articles=1]: 10 ==> [whole milk=1]: 10 <conf:(1)> lift:(3.91) lev:(0) conv:(7.44)
4. [root vegetables=1, whipped/sour cream=1, hygiene articles=1]: 10 ==> [whole milk=1]: 10 <conf:(1)> lift:(3.91) lev:(0) conv:(7.44)
5. [root vegetables=1, whipped/sour cream=1, flour=1]: 17 ==> [whole milk=1]: 17 <conf:(1)> lift:(3.91) lev:(0) conv:(12.66)
6. [root vegetables=1, butter=1, rice=1]: 10 ==> [whole milk=1]: 10 <conf:(1)> lift:(3.91) lev:(0) conv:(7.44)
7. [pip fruit=1, butter=1, hygiene articles=1]: 10 ==> [whole milk=1]: 10 <conf:(1)> lift:(3.91) lev:(0) conv:(7.44)
8. [domestic eggs=1, butter=1, soft cheese=1]: 10 ==> [whole milk=1]: 10 <conf:(1)> lift:(3.91) lev:(0) conv:(7.44)
9. [domestic eggs=1, curd=1, sugar=1]: 10 ==> [whole milk=1]: 10 <conf:(1)> lift:(3.91) lev:(0) conv:(7.44)
10. [domestic eggs=1, napkins=1, cream cheese=1]: 11 ==> [whole milk=1]: 11 <conf:(1)> lift:(3.91) lev:(0) conv:(8.19)
11. [domestic eggs=1, cream cheese=1, sugar=1]: 11 ==> [whole milk=1]: 11 <conf:(1)> lift:(3.91) lev:(0) conv:(8.19)
12. [root vegetables=1, citrus fruit=1, soft cheese=1]: 10 ==> [other vegetables=1]: 10 <conf:(1)> lift:(5.17) lev:(0) conv:(8.07)
13. [pip fruit=1, whipped/sour cream=1, brown bread=1]: 11 ==> [other vegetables=1]: 11 <conf:(1)> lift:(5.17) lev:(0) conv:(8.87)
14. [whole milk=1, rolls/buns=1, soda=1, newspapers=1]: 10 ==> [other vegetables=1]: 10 <conf:(1)> lift:(5.17) lev:(0) conv:(8.07)
15. [other vegetables=1, yogurt=1, root vegetables=1, oil=1]: 14 ==> [whole milk=1]: 14 <conf:(1)> lift:(3.91) lev:(0) conv:(10.42)
16. [whole milk=1, yogurt=1, tropical fruit=1, grapes=1]: 10 ==> [other vegetables=1]: 10 <conf:(1)> lift:(5.17) lev:(0) conv:(8.07)
17. [other vegetables=1, bottled water=1, root vegetables=1, pip fruit=1]: 11 ==> [whole milk=1]: 11 <conf:(1)> lift:(3.91) lev:(0) conv:(8.19)
18. [other vegetables=1, root vegetables=1, butter=1, white bread=1]: 10 ==> [whole milk=1]: 10 <conf:(1)> lift:(3.91) lev:(0) conv:(7.44)
19. [whole milk=1, tropical fruit=1, pip fruit=1, ham=1]: 11 ==> [other vegetables=1]: 11 <conf:(1)> lift:(5.17) lev:(0) conv:(8.87)
20. [other vegetables=1, whipped/sour cream=1, domestic eggs=1, butter=1]: 12 ==> [whole milk=1]: 12 <conf:(1)> lift:(3.91) lev:(0) conv:(8.93)

由上图可以得出, 当支持度为 0.01, 置信度为 0.5 时, 该算法对数据集一可以得出 14 条关联规则, 耗时大约 1-2 秒;



由上图可知，当支持度为 0.001，置信度为 0.9 时，该算法对数据集一可以得出 128 条关联规则，耗时大约 1-2 秒。



由上图可知，当支持度为 0.01，置信度为 0.5 时，该算法对数据集二可得出 78 条关联规则，耗时大约 9 秒。分析与 Apriori 算法不一致的原因可能是在处理数据时出现一条乱码，于是手动将乱码删除，可能导致数据集不一致。

具体的规则详见附件中的输出文件。

六. 结论

在理论上，Apriori 算法比 FPGrowth 算法要消耗更多的时间与空间，因为 Apriori 算法在每一步产生候选项目集时循环产生的组合过多，没有排除不应该参与组合的元素，耗费大量空间；每次计算项集的支持度时，都对数据集中的全部记录进行了一遍扫描比较，耗费大量时间。而 FPGrowth 不使用候选集而且只进行两次数据集扫描。这两种算法与蛮力算法相比都有提升。

在实验中得出的结论也的确如此。处理好数据集后，一开始使用 Weka 进行计算，但是对于给定的两个数据集，Apriori 算法长时间没有得出结果，而 FPGrowth 算法则是只需较短时间便得出结果，后来在 matlab 中手动实现了 Apriori 算法，然而当支持度为 0.001，置信度为 0.5 时，两个数据集仍然均得不出结果，增大支持度到 0.01，减小置信度到 0.5

后才得出结果，并且耗时远大于 FPGrowth 算法的耗时。

七. 参考文献

1. <http://blog.csdn.net/lfdanding/article/details/50755919>
2. <http://www.it165.net/pro/html/201502/33715.html>
3. 数据挖掘概念与技术 原书第三版