



数据分析算法

北京理工大学计算机学院 孙新

2019年1月

4.5 Apriori关联规则算法

- 关联规则挖掘的基本思想
- Apriori 算法的基本过程
- Apriori关联规则算法步骤
- Apriori关联规则算法的特点
- R语言实现Apriori算法示例

- ❖ 利用超市购物篮Groceries数据进行关联规则分析, 使用R语言中的arules包实现Apriori算法
- ❖ 导入包: library(arules); 加载数据源: Groceries数据集

```
> library(arules) #加载 arules 包
> data(Groceries)
> Groceries
transactions in sparse format with
9835 transactions (rows) and
169 items (columns)
```

- ❖ **Groceries数据集**是来自一个现实世界中的超市经营一个月的购物数据，包含了9835次交易，以及169件商品。

[illegible]

- ❖ **数据转换：创建稀疏矩阵，每个Item一列，每一行代表一个transaction。1表示该transaction购买了该item，0表示没有购买。**

```
groceries <- read.transactions("groceries.csv", format="basket", sep=",")
```

ID	Whole milk	...	sausage
1	1	0	1
...	0	1	1
9835	1	0	0

稀疏
矩阵

- ❖ 通过inspect () 函数可以看到超市的交易记录，每次交易的商品名称

```
> inspect(Groceries[1:5])    #通过inspect函数查看Groceries数据集的前5次交易记录
  items
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}
```

通过summary () 函数
可以查看该数据集的一些
基本信息。

```
1 summary(groceries)
2 transactions as itemMatrix in sparse format with
3 9835 rows (elements/itemsets/transactions) and
4 169 columns (items) and a density of 0.02609146
5
6 most frequent items:
7      whole milk other vegetables      rolls/buns      soda
8      2513      1903      1809      1715
9      yogurt      (Other)
10     1372      34055
11
12 element (itemset/transaction) length distribution:
13 sizes
14      1      2      3      4      5      6      7      8      9     10     11     12     13     14     15
15 2159 1643 1299 1005  855  645  545  438  350  246  182  117  78  77  55
16  16  17  18  19  20  21  22  23  24  26  27  28  29  32
17  46  29  14  14   9  11   4   6   1   1   1   1   3   1
18
19      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
20    1.000  2.000  3.000  4.409  6.000 32.000
21
22 includes extended item information - examples:
23      labels
24 1 abrasive cleaner
25 2 artif. sweetener
26 3 baby cosmetics
```

利用itemFrequency () 函数可以查看商品交易比列。也即，取数据集第100到800行，第1列到第3列，计算列代表的三个项目对应的支持度。当然，也可以把支持度itemFrequency排序，查看支持度的最大值，也即取前10个项目对应的支持度。

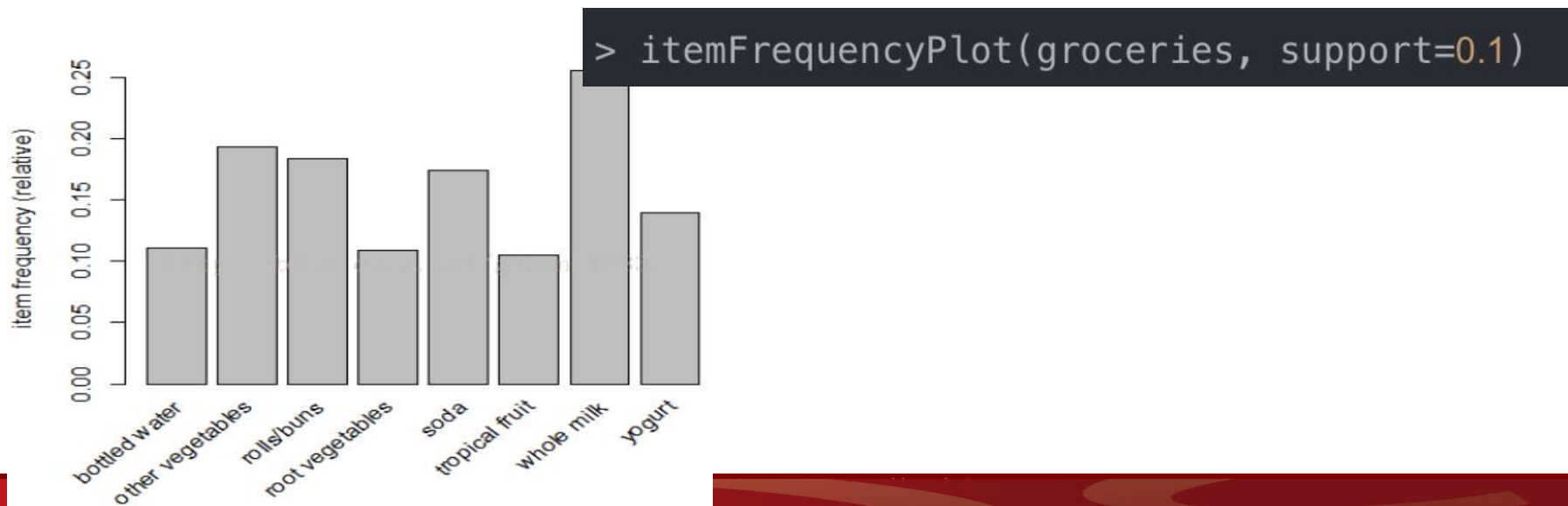
```
1 > itemFrequency(groceries[100:800,1:3])
2 abrasive cleaner artif. sweetener baby cosmetics
3 0.005706134 0.001426534 0.001426534
```

利用sort () 函数可以对支持度进行排序

```
1 > orderedItemFreq <- sort(itemFrequency(groceries), decreasing=T)
2 > orderedItemFreq[1:10]
3 whole milk other vegetables rolls/buns soda yogurt bottled water
4 0.25551601 0.19349263 0.18393493 0.17437722 0.13950178 0.11052364
5 root vegetables tropical fruit shopping bags sausage
6 0.10899847 0.10493137 0.09852567 0.09395018
```

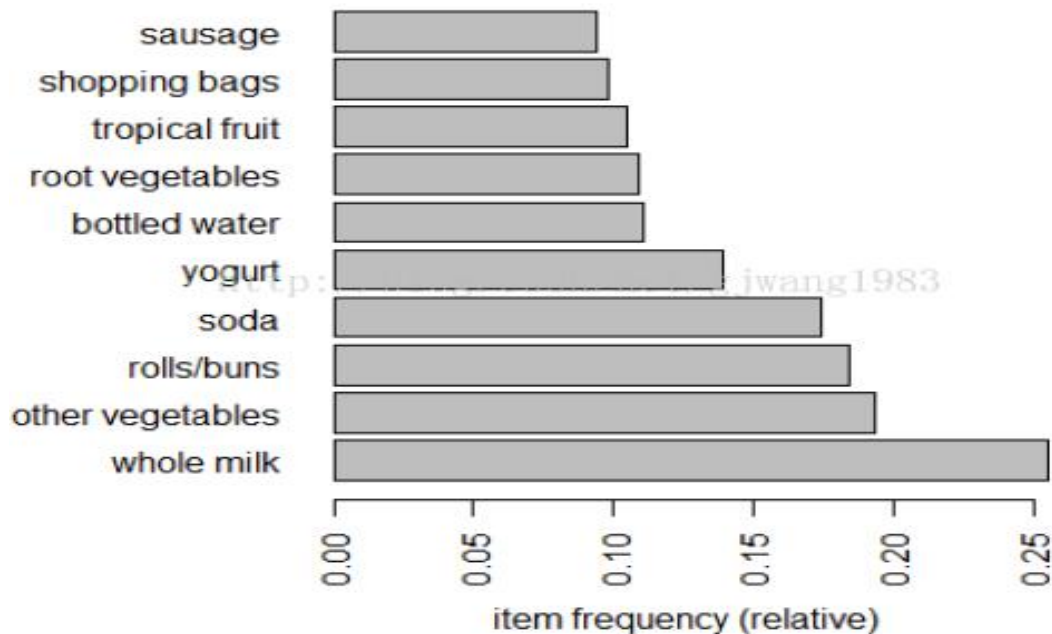

可视化商品的支持度——商品的频率图

为了直观地呈现统计数据，可以使用itemFrequencyPlot()函数生成一个用于描绘所包含的特定商品的交易比例的柱状图。因为包含很多种商品，不可能同时展现出来，因此可以通过support或者topN参数进行排除一部分商品进行展示。support = 0.1 表示支持度至少为0.1。



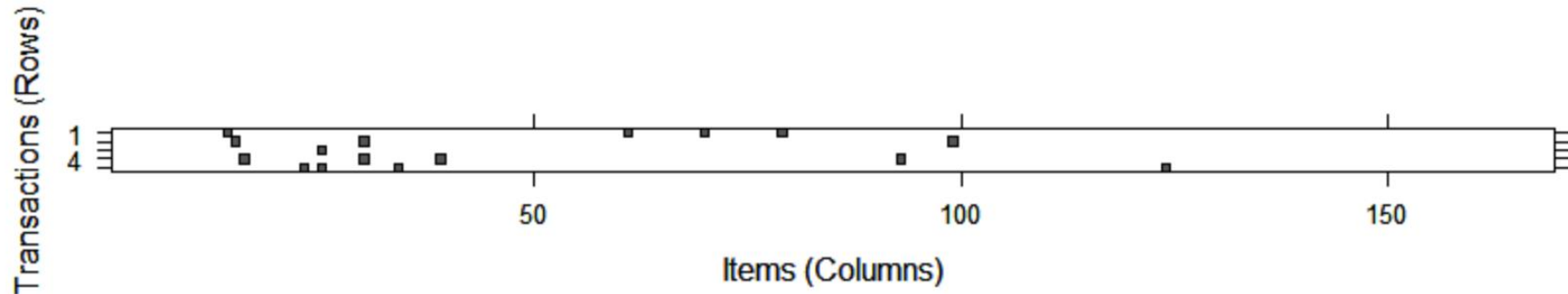
可视化商品的支持度——商品的频率图，topN = 10表示支持度排在前十的商品。

```
> itemFrequencyPlot(groceries, topN=10, horiz=T)
```



可视化交易数据——绘制稀疏矩阵

`image(Groceries[1:5])` # 生成一个5行169列的矩阵，矩阵中填充有黑色的单元表示在此次交易（行）中，该商品（列）被购买了



利用apriori()函数，可以进行规则挖掘。

默认设置 support = 0.1 , confidence = 0.8

```
grocery_rules <- apriori(data=Groceries,parameter=list(support =,confidence =,minlen =))
```

设置支持度和置信度参数来产生合理数量的关联规则时，可能需要进行大量的试验与误差评估。

使用R语言进行规则挖掘，参数设置support = 0.006 , confidence = 0.25

```
1 > groceryrules <- apriori(groceries, parameter = list(support =
2 +                               0.006, confidence = 0.25, minlen = 2))
3
4 Parameter specification:
5 confidence minval smax arem aval originalSupport support minlen maxlen target ext
6      0.25    0.1    1 none FALSE                TRUE  0.006     2    10 rules FALSE
7
8 Algorithmic control:
9 filter tree heap memopt load sort verbose
10    0.1 TRUE TRUE  FALSE TRUE    2    TRUE
11
12 apriori - find association rules with the apriori algorithm
13 version 4.21 (2004.05.09)      (c) 1996-2004 Christian Borgelt
14 set item appearances ... [0 item(s)] done [0.00s].
15 set transactions ... [169 item(s), 9835 transaction(s)] done [0.00s].
16 sorting and recoding items ... [109 item(s)] done [0.00s].
17 creating transaction tree ... done [0.00s].
18 checking subsets of size 1 2 3 4 done [0.01s].
19 writing ... [463 rule(s)] done [0.00s].
20 creating S4 object ... done [0.00s].
```

使用summary () 函数查看规则规汇总信息

```
1 > summary(groceryrules)
2 set of 463 rules
3
4 rule length distribution (lhs + rhs): sizes
5   2   3   4
6 150 297  16
7
8   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
9   2.000  2.000   3.000   2.711  3.000   4.000
10
11 summary of quality measures:
12      support      confidence      lift
13   Min.    :0.006101   Min.    :0.2500   Min.    :0.9932
14   1st Qu.:0.007117   1st Qu.:0.2971   1st Qu.:1.6229
15   Median :0.008744   Median :0.3554   Median :1.9332
16   Mean    :0.011539   Mean    :0.3786   Mean    :2.0351
17   3rd Qu.:0.012303   3rd Qu.:0.4495   3rd Qu.:2.3565
18   Max.    :0.074835   Max.    :0.6600   Max.    :3.9565
19
20 mining info:
21      data ntransactions support confidence
22   groceries          9835    0.006      0.25
```

使用inspect () 查看具体的规则

```
1 > inspect(groceryrules[1:5])
2   lhs                rhs                support confidence    lift
3 1 {potted plants} => {whole milk}      0.006914082    0.4000000 1.565460
4 2 {pasta}          => {whole milk}      0.006100661    0.4054054 1.586614
5 3 {herbs}           => {root vegetables} 0.007015760    0.4312500 3.956477
6 4 {herbs}           => {other vegetables} 0.007727504    0.4750000 2.454874
7 5 {herbs}           => {whole milk}      0.007727504    0.4750000 1.858983
```

使用sort () 对关联规则集合排序

```
1 > ordered_groceryrules <- sort(groceryrules, by="lift")
2 > inspect(ordered_groceryrules[1:5])
3   lhs                rhs                support confidence    lift
4 1 {herbs}              => {root vegetables} 0.007015760 0.4312500 3.956477
5 2 {berries}            => {whipped/sour cream} 0.009049314 0.2721713 3.796886
6 3 {other vegetables,
7   tropical fruit,
8   whole milk}          => {root vegetables} 0.007015760 0.4107143 3.768074
9 4 {beef,
10  other vegetables} => {root vegetables} 0.007930859 0.4020619 3.688692
11 5 {other vegetables,
12  tropical fruit}      => {pip fruit} 0.009456024 0.2634561 3.482649
```


使用itemInfo () 进行查看summary () 函数结果中level1和level2字段的详细信息

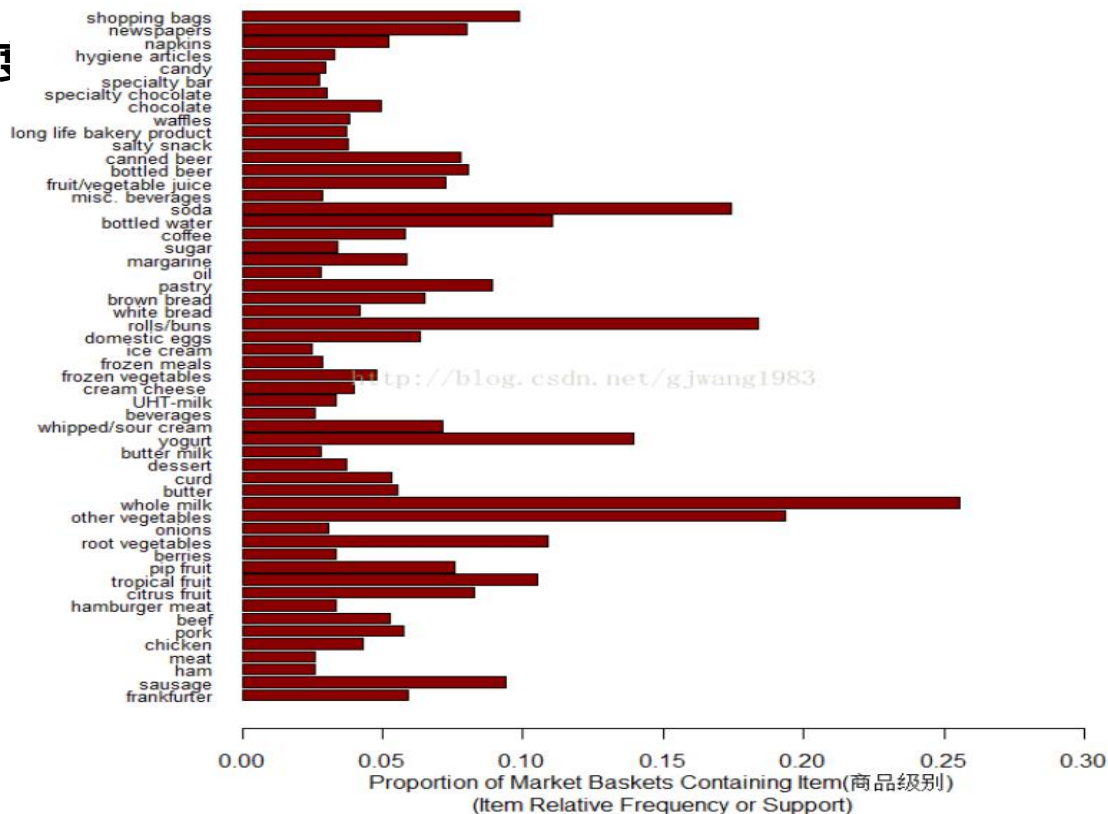
```
> print(levels(itemInfo(Groceries)[["level1"]]))
```

```
[1] "canned food"      "detergent"      "drinks"          "fresh products"  "fruit and vegetables" "meet and sausage"  "non-food"
[8] "perfumery"        "processed food"  "snacks and candies"
```

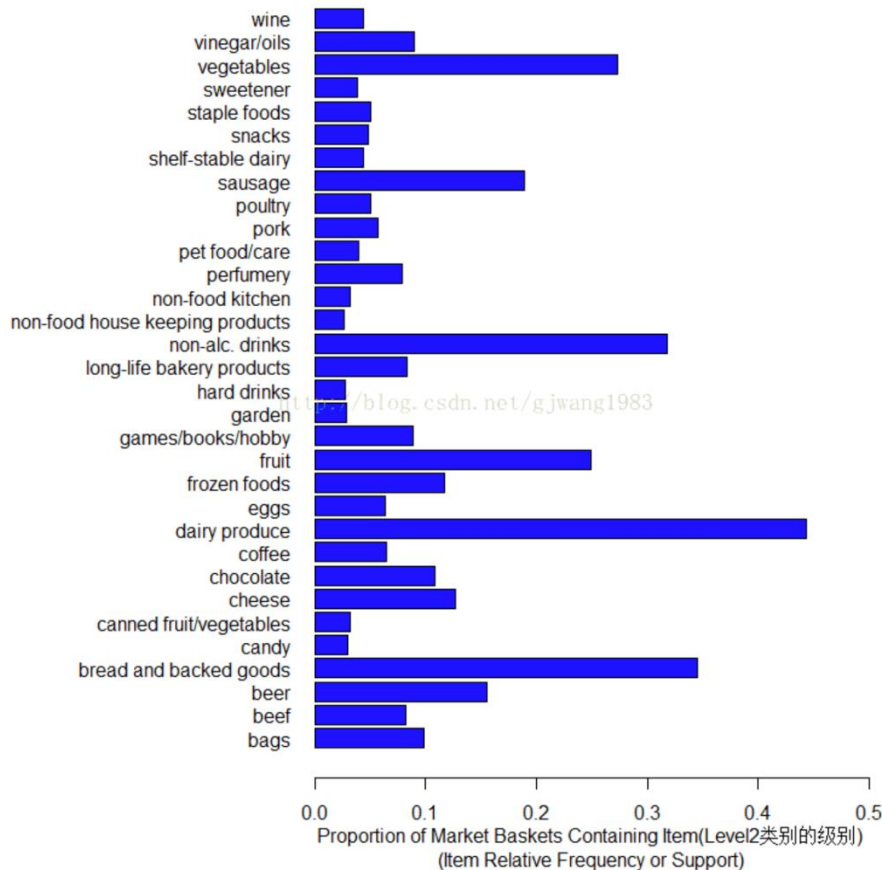
```
> print(levels(itemInfo(Groceries)[["level2"]]))
```

```
[1] "baby food"      "bags"           "bakery improver"  "bathroom cleaner"
[5] "beef"           "beer"           "bread and backed goods" "candy"
[9] "canned fish"    "canned fruit/vegetables" "cheese"           "chewing gum"
[13] "chocolate"     "cleaner"        "coffee"          "condiments"
[17] "cosmetics"      "dairy produce"  "delicatessen"     "dental care"
[21] "detergent/softener" "eggs"          "fish"            "frozen foods"
[25] "fruit"          "games/books/hobby" "garden"          "hair care"
[29] "hard drinks"    "health food"    "jam/sweet spreads" "long-life bakery products"
[33] "meat spreads"   "non-alc. drinks" "non-food house keeping products" "non-food kitchen"
[37] "packaged fruit/vegetables" "perfumery"      "personal hygiene"  "pet food/care"
[41] "pork"           "poultry"        "pudding powder"    "sausage"
[45] "seasonal products" "shelf-stable dairy" "snacks"            "soap"
[49] "soups/sauces"    "staple foods"   "sweetener"         "tea/cocoa drinks"
[53] "vegetables"     "vinegar/oils"   "wine"
```

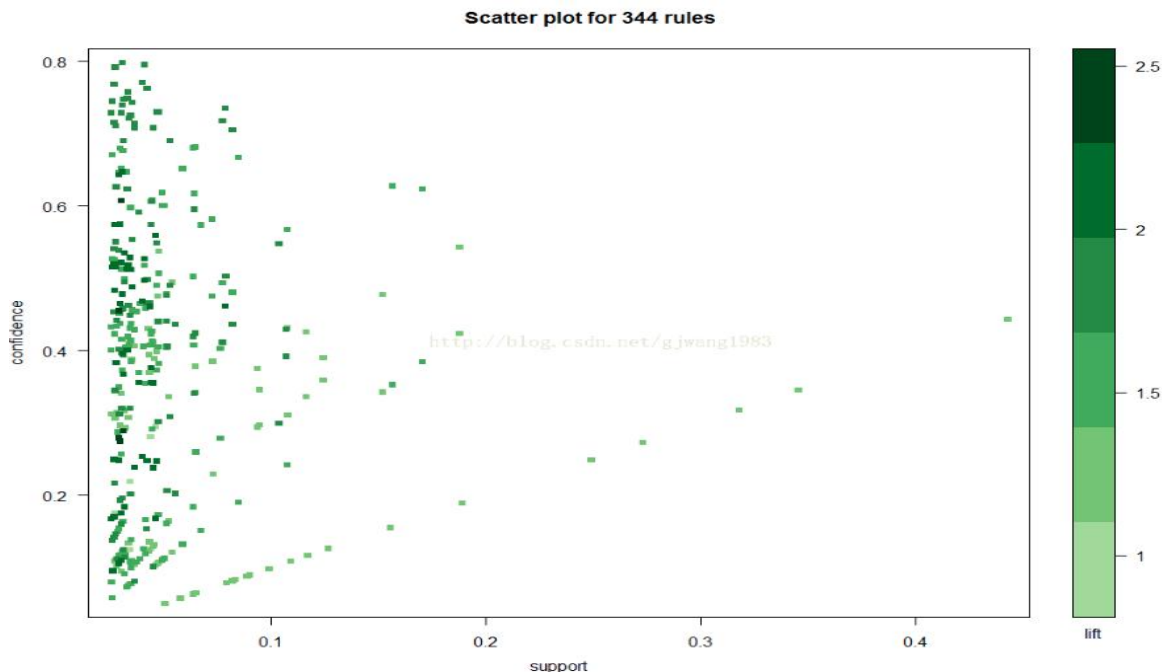
aggregate前的商品 (item) 支持度



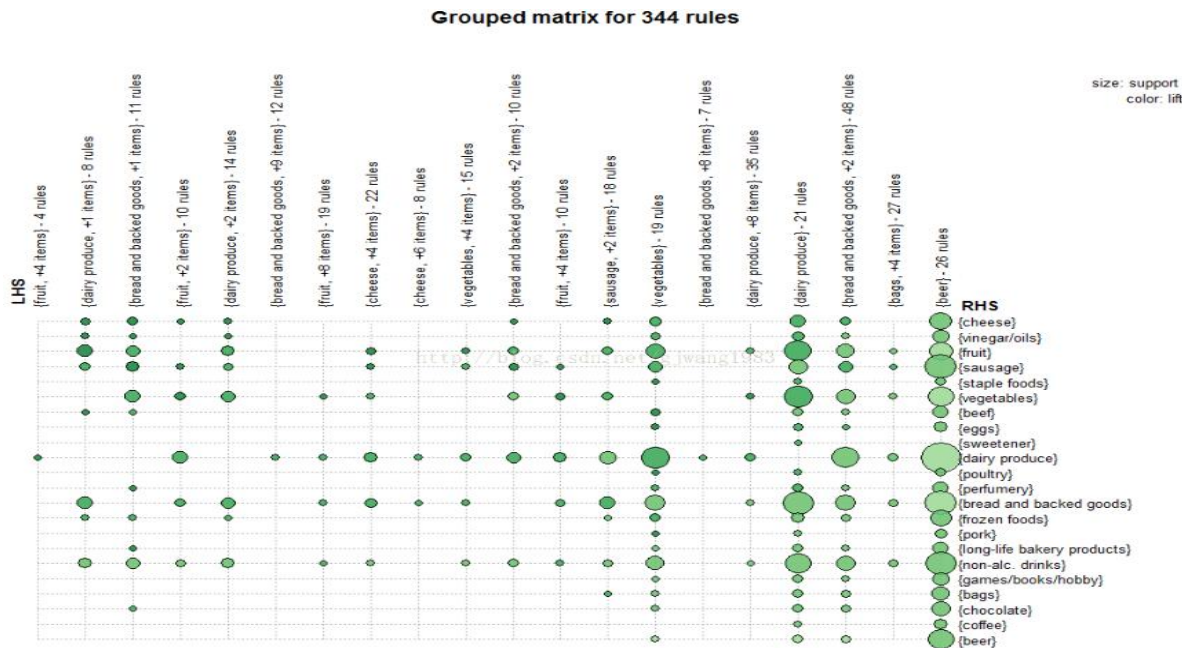
aggregate后类别（category）支持度图：



这幅散点图表示了规则的分布图，横轴支持度，纵轴置信度，颜色深浅为提升度
大部分规则的support在0.1以内，Confidence在0-0.8内。



最后，给出grouped图以有向网状图的形式展示关联规则，图中横坐标为规则前项，纵坐标为规则后项，圆圈表示关联规则，圆圈大小表示支持度大小；颜色深浅代表规则提升度的高低。



谢 谢

Thank you for your attention!