

Oakland Crime频繁模式与关联规则挖掘

目录

- 1. 数据集简介
- 2. 数据集处理与数据项选择
- 3. 数据集频繁模式抽取
- 4. 关联规则计算及评价
- 5. 挖掘结果分析与可视化
 - 5.2 挖掘结果分析
 - 5.1 结果可视化

1. 数据集简介

数据集中包含了2011到2016年奥克兰地区的犯罪记录，每条记录中有11个属性，分别为：

- Agency 机构名称，只有缺失或者“OP”
- Create Time 案件创建时间
- Location 案件发生地点
- Area ID 案件发生区域ID
- Beat 案件发生的巡逻区
- Priority 案件优先级，分别为0, 1, 2
- Incident Type ID 案件类型ID
- Incident Type Description 案件类型
- Event Number 案件标号，唯一值
- Closed Time 案件结束时间
- Zip Codes 案件地区邮政编码

2. 数据集处理与数据项选择

为了挖掘有意义的信息，对此需要对进行挖掘的项(也就是列)进行选择，并将原数据形式转化成利于挖掘的数据格式。鉴于项之间的包含、重叠关系以及项值内容，对项进行如下处理：

为了提取有价值的关联规则，本次实验剔除了无意义的项与唯一项，最终选取的数据属性为：

- Incident Type ID 与 Incident Type Description 二者互相对应，选其一 Incident Type Description 进行实验；
- Area, Beat, Location 均表示案件发生地点，由于 Location 和 Area 属性粒度过细或过粗，这里选取 Beat 属性进行实验；
- Start Time 和 closed Time 属性范围过小，故二者皆取月份作为其特征表示，并根据二者差值新增 Spent Time 项，并根据差值时间不同分为多个类别。

筛选的属性有：

- Agency 机构名称，由于只有缺失值或者OP，剔除
- Event Number 案件标号为一，且不具有规律性，剔除

最终处理后的数据由6项组成，分别为 Incident Type Description、Beat、Start Time、Closed Time 和 Spent time，处理过程如下。

```
1 # 读取数据集
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import datetime
6 %matplotlib inline
7
8 path = "./dataset/oakland_crime_statistics/records-for-{}.csv"
9 data = pd.concat([pd.read_csv(path.format(year)) for year in [2011,2013,2015,2016]])
10 data2 = pd.concat([pd.read_csv(path.format(year)) for year in [2012,2014]])
11 data = pd.concat([data,data2],ignore_index=True)
12 data.info()
```

```
1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 1046388 entries, 0 to 1046387
3 Data columns (total 12 columns):
4 #   Column           Non-Null Count  Dtype  
5 ---  --  
6 0   Agency            1046384 non-null   object 
7 1   Create Time       1046384 non-null   object 
8 2   Location          671477 non-null   object 
9 3   Area Id           864023 non-null   object 
10 4   Beat              1040583 non-null   object 
11 5   Priority          1046384 non-null   float64
12 6   Incident Type Id 1046384 non-null   object 
13 7   Incident Type Description 1045996 non-null   object 
14 8   Event Number      1046384 non-null   object 
15 9   Closed Time       1046359 non-null   object 
16 10  Location 1        374799 non-null   object 
17 11  Zip Codes         352 non-null     float64
18 dtypes: float64(2), object(10)
```

```

1 def str_to_datetime(s):
2     result = s.split('T')
3     date,time = s.split('T')
4     date = date.split('-')
5     time = time[:-4].split(':')
6     date = [int(x) for x in date]
7     time = [int(x) for x in time]
8     return datetime.datetime(date[0],date[1],date[2],time[0],time[1],time[2])
9
10
11 def hierarchy(time):
12     conditions = lambda x: {
13         x < 10: 1,
14         10 <= x < 30: 2,
15         30 <= x < 60: 3,
16         60 <= x < 3*60: 4,
17         3*60 <= x < 6*60: 5,
18         6*60 <= x < 12*60: 6,
19         12*60 <= x < 24*60: 7,
20         x >= 24*60: 8,
21     }
22     return conditions(time)[True]
23
24
25 def time_minus(col1,col2):
26     start = col1.values;
27     end = col2.values;
28
29     ans = []
30     for s,e in zip(start,end):
31         if s=='nan' : e = s
32         elif e == 'nan' : s = e
33         time = (str_to_datetime(e)-str_to_datetime(s)).seconds//60
34         ans.append(hierarchy(time))
35     return ans
36
37 def get_month(col):
38     ans = []
39     for i in col.values:
40         date = str_to_datetime(i)
41         ans.append(date.month)
42     return ans
43
44 def preprocess_data(data):
45     data = data.drop(['Agency','Incident Type Id','Location','Event Number','Area Id','Location 1','Zip Codes'],axis = 1)
46     data['Spent Time'] = time_minus(data['Create Time'],data['Closed Time'])
47     data['Create Time'] = get_month(data['Create Time'])
48     data['Closed Time'] = get_month(data['Closed Time'])
49     return data
50
51 data['Create Time'] = data['Create Time'].astype(str)
52 data['Closed Time'] = data['Closed Time'].astype(str)
53 data = data.dropna(subset=['Closed Time'])
54 data = data.dropna(subset=['Create Time'])
55 data = data[(data['Create Time']!='nan')&(data['Closed Time']!='nan')]
56
57 data = preprocess_data(data)
58 data.head(5)

```

```

1 .dataframe tbody tr th {
2     vertical-align: top;
3 }
4
5 .dataframe thead th {
6     text-align: right;
7 }

```

	Create Time	Beat	Priority	Incident Type Description	Closed Time	Spent Time
0	1	06X	1.0	POSSIBLE DEAD PERSON	1	2
1	1	07X	1.0	415 GUNSHOTS	1	4
2	1	10Y	2.0	415 GUNSHOTS	1	1
3	1	21Y	2.0	415 GUNSHOTS	1	1
4	1	20X	1.0	415 GUNSHOTS	1	3

使用mlxtend库进行数据挖掘，需要将数据进一步处理成其对应格式。

```
1 from mlxtend.frequent_patterns import apriori
2 from mlxtend.preprocessing import TransactionEncoder
3
4 def deal(data):
5     return data.to_list()
6
7 # 增加属性前缀以增加区分度
8 def add_prefix(data):
9     data['Create Time'] = ['s_t' + str(x) for x in data['Create Time'].values]
10    data['Closed Time'] = ['e_t' + str(x) for x in data['Closed Time'].values]
11    data['Spent Time'] = ['l_t' + str(x) for x in data['Spent Time'].values]
12    data['Priority'] = ['p' + str(x) for x in data['Priority'].values]
13    return data
14
15 data = add_prefix(data)
16
17 data['Create Time'] = data['Create Time'].astype(str)
18 data['Closed Time'] = data['Closed Time'].astype(str)
19 data['Spent Time'] = data['Spent Time'].astype(str)
20 data['Priority'] = data['Priority'].astype(str)
21 data['Beat'] = data['Beat'].astype(str)
22 data['Incident Type Description'] = data['Incident Type Description'].astype(str)
23 # 将数据转换为列表，并用TransactionEncoder编码
24 data_list = data.apply(deal, axis=1).tolist()
25 te = TransactionEncoder()
26 data_list_tf = te.fit_transform(data_list)
27 data = pd.DataFrame(data_list_tf, columns=te.columns_)
28
29 data.head(5)
```

	01X	02X	02Y	03X	03Y	04X	05X	05Y	06X	07X	...	s_t11	s_t12	s_t2	s_t3	s_t4	s_t5	s_t
0	False	True	False	...	False	False	False	False	False	False	Fals							
1	False	True	...	False	False	False	False	False	False	Fals								
2	False	...	False	False	False	False	False	False	Fals									
3	False	...	False	False	False	False	False	False	Fals									
4	False	...	False	False	False	False	False	False	Fals									

5 rows × 381 columns

3. 数据集频繁模式抽取

使用Apriori算法计算频繁模式，最小支持度阈值取0.05

```
1 min_support = 0.05
2 frequent_items = apriori(data, min_support=min_support, use_colnames=True, max_len=4).sort_values(by='support', ascending=False)
3 frequent_items.head(100)
```

	support	itemsets
22	0.778700	(p2.0)
18	0.312342	(l_t4)
66	0.234494	(l_t4, p2.0)
21	0.221277	(p1.0)
17	0.172000	(l_t3)
...
42	0.054734	(p2.0, e_t11)
83	0.054459	(s_t11, p2.0, e_t11)
68	0.053733	(p2.0, l_t6)
0	0.052501	(911 HANG-UP)
35	0.052464	(p2.0, 911 HANG-UP)

93 rows × 2 columns

4. 关联规则导出与评价

这里使用的评价方法有教材中的 Lift, Allconf, cosine, Jaccard, Maxconf 以及 Kulczynski , 导出规则表按 Lift 进行降序排序。

```
1 # 使用不同的评价指标进行计算
2 import math
3 from mlxtend.frequent_patterns import association_rules
4
5 def allconf(item):
6     return item.support/max(item['antecedent support'],item['consequent support'])
7 def cosine(item):
8     return item.support/math.sqrt(item['antecedent support']*item['consequent support'])
9 def Jaccard(item):
10    return item.support/(item['antecedent support']+item['consequent support']-item.support)
11 def maxconf(item):
12    return max(item.support/item['antecedent support'],item.support/item['consequent support'])
13 def Kulczynski(item):
14    return 0.5*(item.support/item['antecedent support']+item.support/item['consequent support'])
15
16 def metrics(r,f):
17     ans = []
18     for i in range(r.shape[0]):
19         item = r.iloc[i]
20         ans.append(f(item))
21     return ans
22 rules = association_rules(frequent_items, metric='lift')
23 rules = rules.sort_values(by=['lift'], ascending=False).reset_index(drop=True)
24 rules = rules.drop(['leverage', 'conviction'], axis = 1)
25 rules['cosine'] = metrics(rules,cosine)
26 rules['Jaccard'] = metrics(rules,Jaccard)
27 rules['Allconf'] = metrics(rules,allconf)
28 rules['Maxconf'] = metrics(rules,maxconf)
29 rules['Kulczynski'] = metrics(rules,Kulczynski)
30 rules
```

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	cosine	Jaccard	Allconf	Maxc
0	(s_t11)	(p2.0, e_t11)	0.070393	0.054734	0.054459	0.773650	14.134835	0.877367	0.770648	0.773650	0.994
1	(p2.0, e_t11)	(s_t11)	0.054734	0.070393	0.054459	0.994989	14.134835	0.877367	0.770648	0.773650	0.994
2	(e_t11)	(s_t11, p2.0)	0.070409	0.054744	0.054459	0.773472	14.128860	0.877182	0.770356	0.773472	0.994
3	(s_t11, p2.0)	(e_t11)	0.054744	0.070409	0.054459	0.994798	14.128860	0.877182	0.770356	0.773472	0.994
4	(s_t11)	(e_t11)	0.070393	0.070409	0.070020	0.994705	14.127545	0.994590	0.989239	0.994476	0.994
...
159	(p2.0)	(l_t4)	0.778700	0.312342	0.234494	0.301135	0.964120	0.475479	0.273766	0.301135	0.750
160	(p2.0)	(l_t3)	0.778700	0.172000	0.128378	0.164861	0.958495	0.350784	0.156116	0.164861	0.746
161	(l_t3)	(p2.0)	0.172000	0.778700	0.128378	0.746380	0.958495	0.350784	0.156116	0.164861	0.746
162	(p2.0)	(l_t2)	0.778700	0.141419	0.104859	0.134659	0.952198	0.315985	0.128620	0.134659	0.741
163	(l_t2)	(p2.0)	0.141419	0.778700	0.104859	0.741477	0.952198	0.315985	0.128620	0.134659	0.741

164 rows × 12 columns

5. 挖掘结果分析与可视化

对按 Lift 评分前10条规则进行分析：

```
1 rules = rules.sort_values(by=['lift'], ascending=False).reset_index(drop=True)
2 rules.head(10)
```

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	cosine	Jaccard	Allconf	Maxcon
0	(s_t11)	(p2.0, e_t11)	0.070393	0.054734	0.054459	0.773650	14.134835	0.877367	0.770648	0.773650	0.994981
1	(p2.0, e_t11)	(s_t11)	0.054734	0.070393	0.054459	0.994989	14.134835	0.877367	0.770648	0.773650	0.994981
2	(e_t11)	(s_t11, p2.0)	0.070409	0.054744	0.054459	0.773472	14.128860	0.877182	0.770356	0.773472	0.994791
3	(s_t11, p2.0)	(e_t11)	0.054744	0.070409	0.054459	0.994798	14.128860	0.877182	0.770356	0.773472	0.994791
4	(s_t11)	(e_t11)	0.070393	0.070409	0.070020	0.994705	14.127545	0.994590	0.989239	0.994476	0.994705
5	(e_t11)	(s_t11)	0.070409	0.070393	0.070020	0.994476	14.127545	0.994590	0.989239	0.994476	0.994705
6	(s_t12)	(p2.0, e_t12)	0.073343	0.057119	0.056836	0.774937	13.567059	0.878123	0.771960	0.774937	0.995041
7	(p2.0, e_t12)	(s_t12)	0.057119	0.073343	0.056836	0.995047	13.567059	0.878123	0.771960	0.774937	0.995041
8	(p2.0, s_t12)	(e_t12)	0.057163	0.073297	0.056836	0.994282	13.565110	0.878059	0.771980	0.775422	0.994282
9	(e_t12)	(p2.0, s_t12)	0.073297	0.057163	0.056836	0.775422	13.565110	0.878059	0.771980	0.775422	0.994282

- 所给规则的lift都是远大于1的，说明关联项之间是正相关的
- 0-3项关联规则 表明了优先级为2的案件通常在一个月内(11月)就会解决，优先级为2的案子大多都是小案件，解决速度较快。
- 规则4, 5说明在11月发生的案件可以得到很快的解决。
- 规则6, 9说明在12月优先级为2的案子也会得到很好的解决

5.2 关联规则可视化

x轴和y轴分别代表了关联规则的支持度和置信度，其颜色的深浅则反应了 Lift 评分的大小

```

1 import matplotlib.pyplot as plt
2 plt.xlabel('support')
3 plt.ylabel('confidence')
4 for i in range(rules.shape[0]):
5     plt.scatter(rules.support[i],rules.confidence[i],s=20,c='b',alpha=(rules.lift.iloc[i])/(rules.lift.iloc[0])*0.8/(rules.lift.iloc[0]-rules.lift.iloc[-1])+0.3)

```

