

# 数据挖掘第二次作业 q2 报告

1651718 方沛

产生的数据集 `time_sery.csv` 已经放在文件夹中

## 题目要求

特征工程: 对每个商品编号, 设计如下特征信息:

1. 商品编号、品牌编号、4 级品类结构、日期 (标记为 d)、是否工作日、当日销量、前 1 日 (即 d-1 日) 至前 7 日 (即 d-7 日) 当天销量, 计 16 个特征;
2. 该商品对应品牌前 1 日至前 7 日当日销量, 计 7 个特征;
3. 该商品对应 4 级品类前 1 日至前 7 日当日销量, 计 28 个特征
4. 该商品前第 2 周 (即 d-8 日至 d-14 日)、前第 3 周 (即 d-15 日至 d-21 日)、前第 4 周 (即 d-22 日至 d-28 日) 中该周的每日销量平均值、该周的某日最大值、该周的某日最小值, 计 33=9 个特征;
5. 该品牌前第 2 周 (即 d-8 日至 d-14 日)、前第 3 周 (即 d-15 日至 d-21 日)、前第 4 周 (即 d-22 日至 d-28 日) 中该周的销量每日平均值、该周的某日最大值、该周的某日最小值, 计 33=9 个特征;
6. 该 4 级品牌前第 2 周 (即 d-8 日至 d-14 日)、前第 3 周 (即 d-15 日至 d-21 日)、前第 4 周 (即 d-22 日至 d-28 日) 中该周的销量每日平均值、该周的某日最大值、该周的某日最小值, 计 334=36 个特征;

## 代码设计

### 预处理数据

读入数据集, 设置以 pluno 为索引

```
1 dataset = pd.read_csv("trade_new.csv").fillna(0)
2 dataset['pluno'] = dataset['pluno'].astype('str')
3 dataset['sldatetime'] = pd.to_datetime(dataset['sldatetime'])
4 dataset['sldatetime'] = [datetime.datetime.strptime(x, '%Y-%m-%d') for x in dataset['sldatetime']]
5 dataset = dataset.set_index('pluno', drop = False)
```

插入初始的没有时间序列信息的数据

```
1 for pluno in dataset.index.unique():
2     if str(pluno)[:2] not in ['22', '23', '25', '27']:
3         continue
4     # 一条数据
5     pluno_data = dataset.loc[[pluno]]
6     for index, good in pluno_data.iterrows():
7         row = dict()
8         row['pluno'] = good['pluno']
9         row['bndno'] = good['bndno']
10        row['pl_1'] = int(good['pluno'][:2])
11        row['pl_2'] = int(good['pluno'][:3])
```

```

12     row['pl_3'] = int(good['pluno'][:4])
13     row['pl_4'] = int(good['pluno'][:5])
14     row['purchase_date'] = good['sldatetime']
15     row['is_workday'] = is_workday(good['sldatetime'])
16     row['qty'] = float(good['qty'])
17     row['prev_d1'] = 0
18     row['prev_d1'] = 0
19     row['prev_d2'] = 0
20     row['prev_d3'] = 0
21     row['prev_d4'] = 0
22     row['prev_d5'] = 0
23     row['prev_d6'] = 0
24     row['prev_d7'] = 0
25     ...
26     data_values.append(list(row.values()))
27

```

## 时间序列数据聚合

对每个pluno及其所对应的日期进行聚合

```

1  indexs=['pluno', 'purchase_date']
2  '''
3  @description: 为数据集属性指定聚合函数
4  @params:
5      - properties: 一个dataframe的columns
6  @output: 聚合函数列表
7  '''
8  def get_func(properties):
9      funcs = dict()
10     for prop in properties:
11         if prop == 'pluno' or prop == 'purchase_date':
12             continue
13         elif prop == 'qty':
14             funcs[prop] = 'sum'
15         else:
16             funcs[prop] = 'mean'
17     return funcs
18 # 转换格式才能聚合
19 for col in df.columns:
20     if col == 'pluno' or col == 'purchase_date':
21         continue
22     elif col in ['pl_1', 'pl_2', 'pl_3', 'pl_4']:
23         df[col]=df[col].astype("int")
24     else:
25         df[col]=df[col].astype("float")
26
27 funcs = get_func(properties)
28 df=df.groupby(['pluno', 'purchase_date']).agg(funcs)

```

聚合结果

```
In [12]: 1 df=df.groupby(['pluno','purchase_date']).agg(funcs)
```

```
In [13]: 1 df=df.reset_index()
```

```
In [14]: 1 df.head(5)
```

Out[14]:

	pluno	purchase_date	bndno	pl_1	pl_2	pl_3	pl_4	is_workday	qty	prev_d1	...	pl3_4week_min	pl4_2week_avg	pl4_2week_max	pl4_2week_
0	22000005	2016-07-31	0.0	22	220	2200	22000	0.0	0.704	0.0	...	0.0	0.0	0.0	
1	22000008	2016-02-18	0.0	22	220	2200	22000	1.0	0.704	0.0	...	0.0	0.0	0.0	
2	22000009	2016-07-25	0.0	22	220	2200	22000	1.0	0.666	0.0	...	0.0	0.0	0.0	
3	22000009	2016-07-27	0.0	22	220	2200	22000	1.0	1.120	0.0	...	0.0	0.0	0.0	
4	22000010	2016-03-16	0.0	22	220	2200	22000	1.0	1.914	0.0	...	0.0	0.0	0.0	

5 rows × 105 columns

## 数据工程流程

为了方便调试,将题目中要求的特征生成各封装成了函数

以下函数与题目序号一一对应:

```
1 '''
2 @description: 产生每个pluno过去7天的销量
3 @params:
4     - df:目标数据集
5 '''
6 def set_prev_d(df):
7     count =0
8     pluno_indexs = df.index.levels[0]
9     date_indexs= df.index.levels[1]
10    for p in pluno_indexs:
11        count = count+1
12        print("p:{}".format(p,count))
13        subset = df.loc[p]
14        if subset.shape[0] == 1:
15            continue
16        else:
17            subdates = list(reversed(subset.index))
18            for i in range(len(subdates)):
19                now_date = pd.to_datetime(subdates[i])
20                for j in range(i+1,len(subdates)):
21                    prev_date = pd.to_datetime(subdates[j])
22                    #print(prev_date)
23                    interval_days = int((now_date - prev_date).days)
24                    if 0<interval_days<=7:
25                        col = 'prev_d'+ str(interval_days)
26                        df.at[(p,subdates[i]),col]=df.at[(p,subdates[i]),col]+
df.at[(p,subdates[j]),'qty']
27                        if df.at[(p,subdates[j]),'qty']!=0:
28                            print("index:",p,subdates[j],df.at[(p,subdates[j]),'qty'])
29                            print(df.at[(p,subdates[i]),col])
30                            print('-'*30)
31                    else:
```

```

32                 break
33     '''
34     @description: 产生每个品类过去7天的销量
35     @params:
36         - df:目标数据集
37         - lv:具体哪一级品类
38     '''
39     def set_pl_d(df,lv):
40         all_properties = np.unique([pl[:lv+1] for pl in df['pluno']])
41         for prop in all_properties:
42             print("prop:",prop)
43             subset = df.loc[df['pluno'].str.startswith(prop),:]
44             sorted_subset = subset.sort_values(by=['purchase_date'],ascending=[False])
45             plength=len(sorted_subset.index)
46             for i in range(plength):
47                 now_date = pd.to_datetime(sorted_subset.iloc[i]['purchase_date'])
48                 now = sorted_subset.index[i]
49                 for j in range(i+1,plength):
50                     prev_date = pd.to_datetime(sorted_subset.iloc[j]['purchase_date'])
51                     interval_days = (now_date - prev_date).days
52                     if 0<interval_days<=7:
53                         col = 'pl_{}_d{}'.format(lv,interval_days)
54                         prev = sorted_subset.index[j]
55                         df.at[now,col] =df.at[now,col]+df.at[prev,'qty']
56                         if df.at[prev,'qty']!=0:
57                             print(now,col,df.at[now,col])
58                     else:
59                         break
60     '''
61     @description: 产生每个品类过去几周的最大最小和均值
62     @params:
63         - df:目标数据集
64         - wk:具体哪一周
65     '''
66     def set_pluno_week(df,wk):
67         min_day = 7*(wk-1)
68         max_day = min_day+7
69         all_plunos = np.unique(df['pluno'])
70         colmax = 'pl_{}week_max'.format(wk)
71         colavg = 'pl_{}week_avg'.format(wk)
72         colmin = 'pl_{}week_min'.format(wk)
73         #         tt=0
74         for pluno in all_plunos:
75             print("pluno:",pluno)
76             subset = df.loc[df['pluno']==pluno,:]
77             sorted_subset = subset.sort_values(by=['purchase_date'],ascending=[False])
78
79             plength=len(sorted_subset.index)
80             for i in range(plength):
81                 now_date = pd.to_datetime(sorted_subset.iloc[i]['purchase_date'])
82                 #print("now_date:",now_date)
83                 now = sorted_subset.index[i]
84                 count = 0

```

```

85         for j in range(i+min_day,plength):
86             prev_date = pd.to_datetime(sorted_subset.iloc[j]['purchase_date'])
87             #print("prev_date:",prev_date)
88             interval_days = int((now_date - prev_date).days)
89             #print("interval_days:",min_day,interval_days,max_day)
90             if interval_days>max_day:
91                 break
92             elif min_day<interval_days<=max_day:
93                 count = count+1
94                 prev = sorted_subset.index[j]
95                 if df.at[prev, 'qty']>df.at[now,colmax]:
96                     df.at[now,colmax] = df.at[prev, 'qty']
97                     #print("max:",df.at[now,colmax])
98                 elif df.at[prev, 'qty']<df.at[now,colmin]:
99                     df.at[now,colmin] =df.at[prev, 'qty']
100                    #print("min:",df.at[now,colmax])
101                    df.at[now,colavg]=df.at[now,colavg]+df.at[prev, 'qty']
102                    #print("avg:",df.at[now,colmax])
103
104            if not count == 0:
105                df.at[now,colavg] = df.a
106                t[now,colavg]/count
107            #print(df.at[now,colavg])
108
109
110    '''
111    @description: 产生每个品牌过去几周的最大最小和均值,由于所选的商品都没有品牌,直接复制品类
112    @params:
113        - df:目标数据集
114    '''
115    def set_bndno_week(df):
116        df['bnd_2week_avg'] = df['pl_2week_avg']
117        df['bnd_2week_max'] = df['pl_2week_max']
118        df['bnd_2week_min'] = df['pl_2week_min']
119        df['bnd_3week_avg'] = df['pl_3week_avg']
120        df['bnd_3week_max'] = df['pl_3week_max']
121        df['bnd_3week_min'] = df['pl_3week_min']
122        df['bnd_4week_avg'] = df['pl_4week_avg']
123        df['bnd_4week_max'] = df['pl_4week_max']
124        df['bnd_4week_min'] = df['pl_4week_min']
125
126    '''
127    @description: 产生每个品类过去几周的最大最小和均值
128    @params:
129        - df:目标数据集
130        - wk:哪一周
131        - lv:哪一级品类
132    '''
133    def set_smaller_pluno_week(df,wk,lv):
134        min_day = 7*(wk-1)
135        max_day = min_day+7
136        all_plunos = np.unique([pl[:lv+1] for pl in df['pluno']])
137        colmax = 'pl_{}-{}week_max'.format(lv,wk)

```

```

138     colavg = 'pl_{_}week_avg'.format(lv,wk)
139     colmin = 'pl_{_}week_min'.format(lv,wk)
140     for pluno in all_plunos:
141         print(pluno)
142         subset = df.loc[df['pluno'].str.startswith(pluno),:]
143         sorted_subset = subset.sort_values(by=['purchase_date'],ascending=[False])
144
145         for date in reversed(np.unique(sorted_subset['purchase_date'])):
146             #print("date:",date)
147             indexes = sorted_subset[sorted_subset["purchase_date"]==date].index
148             start_date = pd.to_datetime(date)
149             before_dates = [start_date-dateutil.relativedelta.relativedelta(days=x) for x
in range(min_day,max_day)]
150             before_dates = [datetime.datetime.strftime(x,'%Y-%m-%d') for x in
before_dates]
151             calculated_data =
sorted_subset[sorted_subset["purchase_date"].isin(before_dates)]
152             df.loc[indexs,colmax]=np.max(calculated_data['qty'])
153             df.loc[indexs,colavg]=np.mean(calculated_data['qty'])
154             df.loc[indexs,colmin]=np.min(calculated_data['qty'])

```

在完成这些函数后,我们将它们封装入 `my_feature_engineering_pipeline` 中方便调用

```

1  '''
2  @description: 整个特征工程流程
3  @params:
4      - df:目标数据集
5  '''
6  def my_feature_engineering_pipeline(df):
7      print('step 1')
8      df = df.set_index(['pluno', 'purchase_date'])
9      set_prev_d(df)
10     df =df.reset_index()
11     set_bnd_d(df)
12
13     print('step 2')
14
15     set_pl_d(df,1)
16     set_pl_d(df,2)
17     set_pl_d(df,3)
18     set_pl_d(df,4)
19
20     print('step 3')
21
22     set_pluno_week(df,2)
23     set_pluno_week(df,3)
24     set_pluno_week(df,4)
25
26     print('step 4')
27     set_bndno_week(df)
28     set_smaller_pluno_week(df,2,1)
29     set_smaller_pluno_week(df,3,1)
30     set_smaller_pluno_week(df,4,1)

```

```
31 set_smaller_pluno_week(df,2,2)
32 set_smaller_pluno_week(df,3,2)
33 set_smaller_pluno_week(df,4,2)
34 set_smaller_pluno_week(df,2,3)
35 set_smaller_pluno_week(df,3,3)
36 set_smaller_pluno_week(df,4,3)
37 set_smaller_pluno_week(df,2,4)
38 set_smaller_pluno_week(df,3,4)
39 set_smaller_pluno_week(df,4,4)
```