

构建标准化建模代码（上）





课程目录

ONTENTS

1 变量分箱

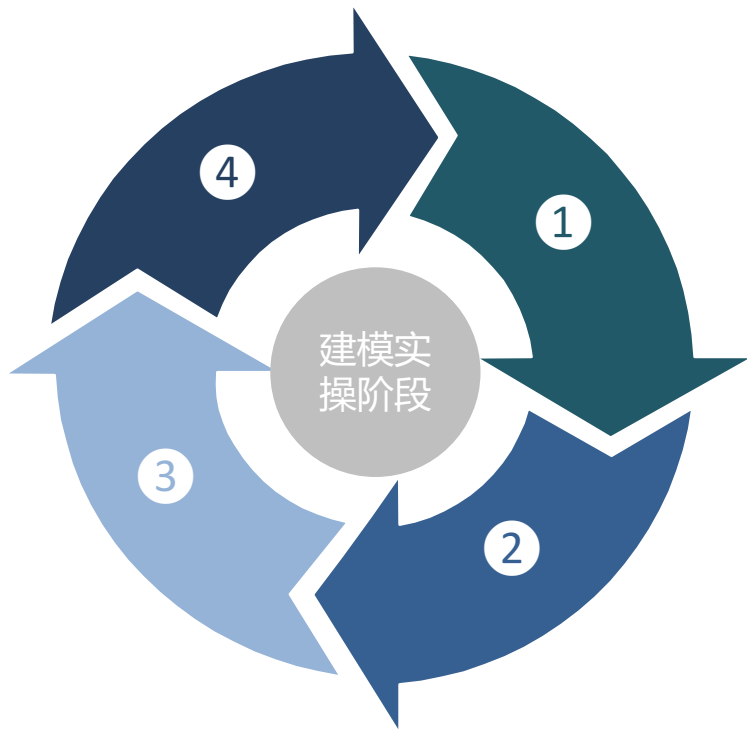
2 变量评价指标计算和woe编码

3 变量筛选

4 模型训练

5 模型评估

6 模型结果标准化输出



① 简单实现

可以针对样本进行建模

② 全流程实现

了解建模全流程的理论，并且可以用code实现，那么就可以独立完成一个评分卡项目

③ 通用化实现

在实现建模全流程code之后，需要思考怎么把代码通用化，应对到不同场景不同项目都可以复用

④ 标准化实现

在实现通用化建模全流程code之后，需要思考如何利用code自动化建模、结果标准化输出最理想的状态就是：输入样本，调一个函数就可以把建模全流程自动走完，自动整合模型结果并标准输出，然后我们需要做的就是坐在旁边等结果

PYTHON3

针对样本量级巨大，也会涉及到Spark

评分卡全流程代码架构

standard_model_tool_20180425_V09

- 01 performance_definition ----- 1 标签定义
 - output_analysis.py ----- 标签定义的分析结果可视化输出
 - vintage_and_roll_rate.py ----- vintage和roll_rate分析
- 02 reject_inference ----- 2 拒绝推断
 - __init__.py
 - reject_inference_method.py ----- 在标签定义之后使用拒绝推断，可以定义多种拒绝推断方法的函数，然后调用
- 03 bin_split ----- 3 变量分箱(分bin)
 - __init__.py
 - best_bin.py ----- 定义变量分bin中的基函数，统计信息函数等功能函数，自动接收切分点输出最终分bin结果
 - ChiMerge_bin.py ----- 卡方分bin方法，输出切分点
 - combine_bin_method.py ----- 合并bin的方法（最大iv、最大gini、最大信息增益等等）
 - cut_bin_method.py ----- 细分bin的方法（最佳ks、最大gini、最大信息增益等等）
 - dynamic_programming_bin.py ----- 动态规划分bin方法，输出切分点
- 04 feature_analyse_and_woe ----- 4 变量评价指标计算和woe编码
 - __init__.py
 - base_feature_analyse_and_woe.py ----- 基于变量分箱的代码，自动线上spark、线下多线程对所有变量进行分箱、评价信息计算、woe转换样本
 - feature_evaluate_method.py ----- 计算变量评价信息的函数，KS、相关系数、单一阈值、MIC、Iv等等
- 05 feature_select ----- 5 变量筛选
 - __init__.py
 - feature_select_method.py ----- 变量筛选功能函数，比如剔除高相关系数函数、逐步回归、随机逻辑回归pvalue检验等等
 - feature_select_rule.py ----- 定义单变量筛选和多变量筛选流程，结合变量评价信息和多变量筛选进行自动变量筛选
- 06 standard_model ----- 6 模型训练
 - __init__.py
 - model_building.py ----- 模型训练code，自动接收筛选后的样本，自动训练模型(切分样本、模型调参、控制训练集测试集稳定性等)，得到模型最终结果
- 07 model_evaluation ----- 7 模型评估
 - __init__.py
 - model_plot.py ----- 模型评估中需要的画图函数，包括ks、roc、PR-F1、lift等等
 - model_result_summary.py ----- 定义模型评估中所有内容的计算函数，并且统一整合
- 08 result_standard_integratation ----- 8 模型结果标准输出
 - __init__.py
 - create_dir.py ----- 定义模型所有结果整合输出的文件夹结构和储存
 - excel_beautify_final.py ----- 整合建模全流程所有结果标准化、美化、自动化输出，可以直接交由评审

1

code demo

变量分箱

① 构建变量分箱函数之前，需要考虑清楚参数

参数名	参数含义
flag_name	标签列名
factor_name	指标列名
data	样本集
bad_name	坏样本个数列名
good_name	好样本个数列名
piece	最大分组数
rate	每组样本最小占比
min_bin_size	每组样本最小数
not_in_list	缺失值列表
cnt_method	切分方法
combine_method	聚合方法
bin_thought	分bin思想
if_need_control	是否需要woe列表进行顺序控制



② 然后将样本实现如下初步统计

Index	var	#Obs	#Good	#Bad	%Bad_Rate
0	317	1	0	1	100.00%
1	329	2	1	1	50.00%
2	332	14	0	14	100.00%
3	333	6	2	4	66.67%
4	334	9	0	9	100.00%
5	335	4	0	4	100.00%
6	336	7	0	7	100.00%
7	338	37	0	37	100.00%
8	339	7	0	7	100.00%
9	340	14	1	13	92.86%
10	341	23	3	20	86.96%
11	342	20	0	20	100.00%
12	345	1	0	1	100.00%
13	346	3	0	3	100.00%
14	348	7	4	3	42.86%
15	350	2	0	2	100.00%
16	351	11	2	9	81.82%
17	352	19	1	18	94.74%
18	353	8	3	5	62.50%
19	354	26	1	25	96.15%
20	355	13	5	8	61.54%
21	356	13	5	8	61.54%
22	357	23	17	6	26.09%
23	358	19	4	15	78.95%
24	359	33	4	29	87.88%
25	360	28	9	19	67.86%
26	361	17	4	13	76.47%
27	362	12	4	8	66.67%
28	363	14	1	13	92.86%



③ 根据不同的分bin思想，获得最佳切分点

(左开右闭)

Index	var	#Obs	#Good	#Bad	%Bad_Rate
0	317	1	0	1	100.00%
1	329	2	1	1	50.00%
2	332	14	0	14	100.00%
3	333	6	2	4	66.67%
4	334	9	0	9	100.00%
5	335	4	0	4	100.00%
6	336	7	0	7	100.00%
7	338	37	0	37	100.00%
8	339	7	0	7	100.00%
9	340	14	1	13	92.86%
10	341	23	3	20	86.96%
11	342	20	0	20	100.00%
12	345	1	0	1	100.00%
13	346	3	0	3	100.00%
14	348	7	4	3	42.86%
15	350	2	0	2	100.00%
16	351	11	2	9	81.82%
17	352	19	1	18	94.74%
18	353	8	3	5	62.50%
19	354	26	1	25	96.15%
20	355	13	5	8	61.54%
21	356	13	5	8	61.54%
22	357	23	17	6	26.09%
23	358	19	4	15	78.95%
24	359	33	4	29	87.88%
25	360	28	9	19	67.86%
26	361	17	4	13	76.47%
27	362	12	4	8	66.67%
28	363	14	1	13	92.86%

④ 根据切分点统计分bin结果

Characteristic	Description	Bin	WOE	IV_Bin	IV_Total	KS	#Obs
RH_out	变量含义	(-inf, 57.5]	0.83	0.05	0.06	0.17	844
		(57.5, 63.0]	0.10	0.00	0.06	0.07	459
		(63.0, 71.3]	0.00	0.00	0.06	0.09	994
		(71.3, 88.5]	-0.03	0.00	0.06	0.07	3244
		(88.5, inf)	-0.15	0.01	0.06	0.05	3156



#Good	#Bad	%Obs	%Bad	%Cumulative_Bad	%Bad_Rate	%Cumulative_Bad_Rate
772	72	9.70%	4.68%	4.68%	8.53%	8.53%
384	75	5.28%	4.88%	9.56%	16.34%	11.28%
818	176	11.43%	11.44%	21.00%	17.71%	14.06%
2657	587	37.30%	38.17%	59.17%	18.09%	16.42%
2528	628	36.29%	40.83%	100.00%	19.90%	17.68%


```
152 @staticmethod
153 def group_by_df(data, flag_name, factor_name, bad_name, good_name, discrete_list=[]):
154     """method's help string.
155
156     Explanation
157     转换数据格式，变为指标每一种取值、每一种指标取值下好样本的个数、坏样本的个数
158
159     Parameters
160
161     data: dataframe
162     | 原始样本
163
164     flag_name: string
165     | 标签名称
166
167     factor_name: string
168     | 指标名称
169
170     bad_name: string
171     | 坏样本个数列名
172
173     good_name: string
174     | 好样本个数列名
175
176     discrete_list: list, default=[]
177     | 指定离散指标的名称list
178     """
179
180     if len(data) == 0:
181         return pandas.DataFrame()
182     regroup1 = data.groupby([factor_name])[flag_name].count()
183     regroup2 = data.groupby([factor_name])[flag_name].sum()
184     data1 = pandas.DataFrame({good_name: regroup1-regroup2, bad_name: regroup2}).reset_index()
185     data1['%Bad_Rate'] = data1[bad_name]/(data1[bad_name]+data1[good_name])
186     if factor_name not in discrete_list:
187         #data1[factor_name] = data1[factor_name].astype(float)
188         data1 = data1.sort_values(by=[factor_name], ascending=True)
189         #data1[factor_name] = list(map(lambda x: Best Bin Cut.get_str(x), data1[factor_name]))
190         data1['Char_Type'] = 'numeric'
191     else:
192         data1 = data1.sort_values(by=['%Bad_Rate'], ascending=True)
```

```
193     data1['Char_Type'] = 'non-numeric'
194 data1 = data1.reset_index(drop=True)
195 return data1
```

```
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
262
263
264
265
266
267
268
269
270
271
272
273
```

③ 计算最佳分组切分点

```
1  -*- coding:utf-8 -*-
2  # creator:
3  import pandas
4  import numpy
5  import math
6  def ChiMerge_Bin(df, good_name='#Good', bad_name='#Bad', confidenceVal=3.841, maxbin=10, minbin=3, minBinPcnt=0):
7      '''变量分箱
8      Parameters
9      -----
10     df: dataframe, 传入一个数据框仅包含一个需要卡方分箱的变量与正负样本标识(响应样本为1, 非响应为0)
11     var: 变量名. var需在df中
12     flag: y_train
13     confidenceVal:  $P(K^2 \geq k) = 0.05$  表示组间相近的概率为5%, k越大, p越小
14     maxbin: 最多箱子数量
15     minbin: 最少箱子数量
16     minBinPcnt: 组内样本最低占比
17     Returns
18     -----
19     result_data: dataframe
20     包含该变量分箱点, woe, iv等信息
21     '''
22     var = df.columns[0]
23     regroup = df[[var, bad_name, good_name]]
24     np_regroup = numpy.array(regroup) # 把数据框转化为numpy (提高运行效率)
25
26     # 处理连续没有正样本或负样本的区间, 并进行区间的合并 (以免卡方值计算报错)
27     i = 0
28     while (i <= np_regroup.shape[0] - 2):
29         if ((np_regroup[i, 1] == 0 and np_regroup[i + 1, 1] == 0) or (
30             np_regroup[i, 2] == 0 and np_regroup[i + 1, 2] == 0)):
31             np_regroup[i, 1] = np_regroup[i, 1] + np_regroup[i + 1, 1] # 正样本
32             np_regroup[i, 2] = np_regroup[i, 2] + np_regroup[i + 1, 2] # 负样本
33             np_regroup[i, 0] = np_regroup[i + 1, 0]
34             np_regroup = numpy.delete(np_regroup, i + 1, 0)
35             i = i - 1
36         i = i + 1
37     # 先保证最后一组大于5%
38     i = np_regroup.shape[0] - 1
39     total_len = float(sum(np_regroup[:, 1] + np_regroup[:, 2]))
40     while (i > 0):
41         if ((np_regroup[i, 1] + np_regroup[i, 2]) / total_len >= minBinPcnt):
42             break
```

```

43         break
44     if ((np_regroup[i, 1] + np_regroup[i, 2])/total_len < minBinPcnt):
45         np_regroup[i, 1] = np_regroup[i, 1] + np_regroup[i - 1, 1] # 响应样本
46         np_regroup[i, 2] = np_regroup[i, 2] + np_regroup[i - 1, 2] # 非响应样本
47         np_regroup = numpy.delete(np_regroup, i - 1, 0)
48     i = i - 1
49 #计算分组样本占比, 把低于5%的分组合并
50 i = 0
51 total_len = float(sum(np_regroup[:, 1] + np_regroup[:, 2]))
52 while (i <= np_regroup.shape[0] - 2):
53     if ((np_regroup[i, 1] + np_regroup[i, 2])/total_len < minBinPcnt):
54         np_regroup[i, 1] = np_regroup[i, 1] + np_regroup[i + 1, 1] # 响应样本
55         np_regroup[i, 2] = np_regroup[i, 2] + np_regroup[i + 1, 2] # 非响应样本
56         np_regroup[i, 0] = np_regroup[i + 1, 0]
57         np_regroup = numpy.delete(np_regroup, i + 1, 0)
58         i = i - 1
59     i = i + 1
60 # 对相邻两个区间进行卡方值计算
61 chi_table = numpy.array([]) # 创建一个数组保存相邻两个区间的卡方值  $n(ad-bc)^2/(a+b)(c+d)(a+c)(b+d)$ 
62 for i in numpy.arange(np_regroup.shape[0] - 1):
63     a = np_regroup[i, 1]
64     b = np_regroup[i, 2]
65     d = np_regroup[i + 1, 2]
66     c = np_regroup[i + 1, 1]
67     chi = (a + b + c + d) * (a*d - b*c) ** 2 / ((a + b) * (c + d) * (a + c) * (b + d))
68     chi_table = numpy.append(chi_table, chi)
69 print('已完成数据初处理, 正在进行卡方分箱核心操作')
70 # 把卡方值最小的两个区间进行合并 (卡方分箱核心)
71 while (1):
72     if (len(chi_table) >= (minbin - 1)):
73         if (len(chi_table) <= (maxbin - 1) and min(chi_table) >= confidenceVal):
74             break
75         else:
76             chi_min_index = numpy.argwhere(chi_table == min(chi_table))[0] # 找出卡方值最小的位置索引
77             np_regroup[chi_min_index, 1] = np_regroup[chi_min_index, 1] + np_regroup[chi_min_index + 1, 1]
78             np_regroup[chi_min_index, 2] = np_regroup[chi_min_index, 2] + np_regroup[chi_min_index + 1, 2]
79             np_regroup[chi_min_index, 0] = np_regroup[chi_min_index + 1, 0]
80             np_regroup = numpy.delete(np_regroup, chi_min_index + 1, 0) #删除指定行
81             '''
82             计算合并后分组相关联的卡方值并替换
83             1、当chi_min_index位于chi_table最后一个时: 计算合并后分组与前一个分组的卡方值
84             2、当chi_min_index位于chi_table第一个时: 计算合并后分组与后一个分组的卡方值
85             3、否则: 计算合并后分组与前一个分组、后一个分组的卡方值
86             '''
87             if (chi_min_index == np_regroup.shape[0] - 1):

```

```

88     a = np_rerroup[chi_min_index - 1, 1]
89     b = np_rerroup[chi_min_index - 1, 2]
90     c = np_rerroup[chi_min_index, 1]
91     d = np_rerroup[chi_min_index, 2]
92     chi_table[chi_min_index - 1] = (a + b + c + d) * (a*d - b*c) ** 2 / ((a + b) * (c + d) * (a + c) * (b + d)) #计算与前一组的卡方值
93     chi_table = numpy.delete(chi_table, chi_min_index, axis=0) # 删除替换前的卡方值
94     elif chi_min_index == 0:
95         a1 = np_rerroup[chi_min_index, 1]
96         b1 = np_rerroup[chi_min_index, 2]
97         c1 = np_rerroup[chi_min_index + 1, 1]
98         d1 = np_rerroup[chi_min_index + 1, 2]
99         chi_table[chi_min_index] = (a1 + b1 + c1 + d1) * (a1*d1 - b1*c1) ** 2 / ((a1 + b1) * (c1 + d1) * (a1 + c1) * (b1 + d1)) #计算与后一组的卡方值
100        chi_table = numpy.delete(chi_table, chi_min_index + 1, axis=0) #删除被替换组的卡方值
101
102    else:
103        a = np_rerroup[chi_min_index - 1, 1]
104        b = np_rerroup[chi_min_index - 1, 2]
105        c = np_rerroup[chi_min_index, 1]
106        d = np_rerroup[chi_min_index, 2]
107        chi_table[chi_min_index - 1] = (a + b + c + d) * (a*d - b*c) ** 2 / ((a + b) * (c + d) * (a + c) * (b + d)) #计算与前一组的卡方值
108
109        a1 = np_rerroup[chi_min_index, 1]
110        b1 = np_rerroup[chi_min_index, 2]
111        c1 = np_rerroup[chi_min_index + 1, 1]
112        d1 = np_rerroup[chi_min_index + 1, 2]
113        chi_table[chi_min_index] = (a1 + b1 + c1 + d1) * (a1*d1 - b1*c1) ** 2 / ((a1 + b1) * (c1 + d1) * (a1 + c1) * (b1 + d1)) #计算与后一组的卡方值
114
115        chi_table = numpy.delete(chi_table, chi_min_index + 1, axis=0) #删除被替换组的卡方值
116
117    else:
118        break
119    #平滑分割点
120    cutpoint=list(regroup[regroup[var].isin(np_rerroup[:, 0])].index)[0:-1]
121    return cutpoint
122
123
124
125
126
127
128
129
130
131
132

```

4 根据切分点统计分组信息

```
593
594 @staticmethod
595 def important_indicator_calculator(data_df, good_name, bad_name, factor_name, knots_list, na_df, na_zero):
596     """
597     Explanation
598     -----
599     计算最大ks对应的点
600
601     Parameters
602     -----
603     data_df: dataframe
604     | 转换后的数据
605
606     good_name: string
607     | 好样本个数列名
608
609     bad_name: string
610     | 坏样本个数列名
611
612     factor_name: string
613     | 指标列名
614
615     knots_list: list
616     | 最佳分组点集合
617
618     na_df: dict
619     | 指标值单独分为一组的样本
620
621     na_zero: Boolean
622     | 缺失值数量是否小于最小分组数量阈值或者最小分组占比阈值，如果低于阈值，则WOE为0，否则正常计算
623
624     Return
625     -----
626     result_indicator: dataframe
627     | 分bin结果
628     """
629     flag = data_df['Char_Type'].max()
630     temp_df_list = []
631     bin_list = []
632     for i in range(1, len(knots_list)):
633         if i == 1:
```

```

634 temp_df_list.append(data_df.loc[knots_list[i - 1]:knots_list[i]])
635 if flag=='numeric':
636     bin_list.append('(-inf, ' + Best_Bin_Cut.get_str(data_df[factor_name][knots_list[i]]) + ')')
637 else:
638     bin_list.append(list(data_df[factor_name])[knots_list[i - 1]:knots_list[i]+1])
639 else:
640     # if knots_list[i - 1]==knots_list[i]:
641     #     temp_df_list.append(data_df.loc[knots_list[i - 1] :knots_list[i]])
642     # else:
643     #     if knots_list[len(knots_list) - 2]==knots_list[len(knots_list)-1] and i==len(knots_list)-2:
644     #         temp_df_list.append(data_df.loc[knots_list[i - 1] + 1:knots_list[i]-1])
645     #     else:
646     #         temp_df_list.append(data_df.loc[knots_list[i - 1] + 1:knots_list[i]])
647     temp_df_list.append(data_df.loc[knots_list[i - 1] + 1:knots_list[i]])
648     if flag=='numeric':
649         if i == len(knots_list) - 1:
650             bin_list.append('(' + Best_Bin_Cut.get_str(data_df[factor_name][knots_list[i - 1]]) + ', inf)')
651         else:
652             bin_list.append(
653                 '(' + Best_Bin_Cut.get_str(data_df[factor_name][knots_list[i - 1]]) + ', ' + Best_Bin_Cut.get_str(data_df[factor_name][knots_list[i]])
654             )
655         else:
656             bin_list.append(list(data_df[factor_name])[knots_list[i - 1] + 1:knots_list[i]+1])
657 if len(knots_list)==2:
658     bin_list=[(-inf, inf)]
659 bin_type_list=['Order_Bin']*len(bin_list)
660 if len(na_df) != 0:
661     tmp_value=list(na_df.values())
662     total_good = sum(data_df[good_name]) + sum([sum(key[good_name]) for key in tmp_value])
663     total_bad = sum(data_df[bad_name]) + sum([sum(key[bad_name]) for key in tmp_value])
664     temp_df_list.extend(tmp_value)
665     bin_list.extend(na_df.keys())
666     bin_type_list=bin_type_list+['Single_Bin']*len(tmp_value)
667 else:
668     total_good = sum(data_df[good_name])
669     total_bad = sum(data_df[bad_name])
670 good_percent_series = pandas.Series(list(map(lambda x: float(sum(x[good_name])) / total_good, temp_df_list)))
671 bad_percent_series = pandas.Series(list(map(lambda x: float(sum(x[bad_name])) / total_bad, temp_df_list)))
672 woe_list = list(numpy.log(good_percent_series / bad_percent_series))
673 IV_list = list((good_percent_series - bad_percent_series) * numpy.log(good_percent_series / bad_percent_series))
674 tmp_key=list(na_zero.keys())
675 for zero in na_zero:
676     if na_zero[zero]:
677         tmp_num=len(tmp_key)-tmp_key.index(zero)
678         woe_list[-tmp_num]=0
679         IV_list[-tmp_num]=0

```

```

680 # if na_zero:
681     #     woe_list=woe_list[0:-1]+[0]
682     #     IV_list=IV_list[0:-1]+[0]
683     IV_total=sum(IV_list)
684     KS_list=list(map(lambda x:max(abs(numpy.cumsum(x[good_name]) / float(sum(x[good_name]))-numpy.cumsum(x[bad_name]) / float(sum(x[bad_name])))),temp_df_list))
685     ks_total=max(abs(numpy.cumsum(data_df[good_name]) / float(sum(data_df[good_name]))-numpy.cumsum(data_df[bad_name]) / float(sum(data_df[bad_name]))))
686     good_list = list(map(lambda x: sum(x[good_name]), temp_df_list))
687     bad_list = list(map(lambda x: sum(x[bad_name]), temp_df_list))
688     total_list = list(map(lambda x: sum(x[good_name])+sum(x[bad_name]), temp_df_list))
689     total_bad_rate_list = list(map(lambda x: float(sum(x[bad_name])) / total_bad, temp_df_list))
690     cumsum_total_bad_rate_list=numpy.cumsum(total_bad_rate_list)
691     bin_bad_rate_list = list(
692         map(lambda x: float(sum(x[bad_name])) / (sum(x[good_name]) + sum(x[bad_name])), temp_df_list))
693     bin_rate_list=list(map(lambda x: float(sum(x[good_name]) + sum(x[bad_name])) / (total_good+total_bad), temp_df_list))
694     non_na_indicator = pandas.DataFrame({'Bin': bin_list,
695                                         'Bin_Type':bin_type_list,
696                                         'WOE': woe_list,
697                                         'IV_Bin': IV_list,
698                                         'IV_Total':IV_total,
699                                         'KS':KS_list,
700                                         'KS_Total':ks_total,
701                                         '#Obs':total_list,
702                                         '#Good':good_list,
703                                         '#Bad': bad_list,
704                                         '%Obs': bin_rate_list,
705                                         '%Bad':total_bad_rate_list,
706                                         '%Cumulative_Bad': cumsum_total_bad_rate_list,
707                                         '%Bad_Rate': bin_bad_rate_list,
708                                         'Char_Type': flag})
709     non_na_indicator['%Cumulative_Bad_Rate'] = numpy.cumsum(non_na_indicator[bad_name])/numpy.cumsum(non_na_indicator['#Obs'])
710     result_indicator = non_na_indicator.reset_index(drop=True)
711     result_indicator['Characteristic'] = factor_name
712     cop_bad_rate=list(result_indicator[result_indicator['Bin_Type']=='Order_Bin']['%Bad_Rate'])
713     result_indicator['Char_Direction'] = Best_Bin_Cut.trend_analyse(cop_bad_rate)
714     columns_order=['Characteristic','Char_Type','Bin','Char_Direction','Bin_Type','WOE','IV_Bin','IV_Total','KS','KS_Total','#Obs','#Good','#Bad','%Obs','%Bad']
715     result_indicator=result_indicator[columns_order]
716     return result_indicator
717
718
719
720
721
722
723

```


2

code demo

变量评价指标 计算和woe编码



① 变量评价指标函数

```
44 class Feature_Evaluate_Method(object):
45
46     def __init__(self):
47         pass
48
49     @staticmethod
50     def single_threshold(data, factor_name):
51         """
52         Explanation
53         -----
54         计算单一阈值
55
56         Remark
57         (1) 适用于连续性指标、离散有序型指标、离散无序型指标
58         (2) 适用于粗分类后的新指标(哑变量或者woe证据权重)
59         (3) 缺失值也参与计算
60         """
61         percent = data[factor_name].value_counts(normalize=True, dropna=False)
62         return percent.max()
63
64     @staticmethod
65     def correlation(data, flag_name, factor_name, corr_method):
66         """
67         Explanation
68         -----
69         计算相关系数
70
71         Remark
72         (1) 适用于连续性指标、离散有序型指标
73         (2) 适用于粗分类后的新指标(哑变量或者woe证据权重)
74         (3) 如果指标中存在缺失值，剔除缺失值计算
75         """
76         try:
77             data[factor_name] = data[factor_name].astype(float)
78             return data[[flag_name, factor_name]].corr(method=corr_method)[flag_name][factor_name]
79         except:
80             return None
81
82     @staticmethod
83     def KS(data, flag_name, factor_name):
84         """
```

Explanation

计算KS,以及KS对应的指标取值

Remark

- (1) 适用于连续性指标、离散有序型指标、离散无序型指标
- (2) 适用于粗分类后的新指标(哑变量或者woe证据权重)
- (3) 如果指标中存在缺失值,剔除缺失值计算

```
data1 = data[[flag_name, factor_name]]
if len(data1)==0:
    return None,None
data1=data1[data1[factor_name].notnull()]
data1=Best_Bin_Cut.group_by_df(data1,flag_name,factor_name,'bad','good')
if len(data1)==0:
    return 0,None
default_CDF = numpy.cumsum(data1['bad']) / sum(data1['bad'])
undefault_CDF = numpy.cumsum(data1['good']) / sum(data1['good'])
ks_CDF=abs(default_CDF - undefault_CDF)
ks = max(ks_CDF)
if numpy.isnan(ks):
    return None,None
index=list(filter(lambda x:ks_CDF.loc[x]==ks,ks_CDF.index))[0]
value=data1[factor_name][index]
return ks,value
```

@staticmethod

```
def IV(data,flag_name,factor_name):
```

Explanation

Remark

- (1) 适用于连续性指标、离散有序型指标、离散无序型指标
- (2) 适用于粗分类后的新指标(哑变量或者woe证据权重)
- (3) 如果指标中存在缺失值,剔除缺失值计算

```
data1=data[[flag_name,factor_name]]
data1=data1[data1[factor_name].notnull()]
data1=Best_Bin_Cut.group_by_df(data1,flag_name,factor_name,'bad','good')
good=float(sum(data1['good']))
bad=float(sum(data1['bad']))
data1['iv']=list(map(lambda x,y:(x/good-y/bad)*numpy.log(x*bad/(good*y)) if good*y!=0 else numpy.inf,data1['good'],data1['bad']))
return data1['iv'].sum()
```

轻量级

样本描述：样本数20000；变量数6000

计算方法：[线下python多进程计算](#)

代码效率：完成所有变量的分箱只需要1.5min
~ 2min；分箱的同时再计算变量评价指标、
woe转换，只需要4min ~ 5min

Ps：线下即使不使用多进程，通过优化分箱
代码，循环所有变量可以使得速度达到3min

重量级

样本描述：样本数100万；变量数40000

计算方法：[线上spark分布式计算](#)

代码效率：完成所有变量的分箱只需要15min
~ 25min；分箱的同时再计算变量评价指标、
woe转换，只需要25min ~ 35min



在数据量达到100万，变量数达到40000时，针对单个变量进行分箱，python尚能支持，**但是对40000变量全部分箱，python根本无法运行起来**，必须结合spark分布式计算，但是了解spark的同学都知道，**spark的map过程，每次传入的对象是一行**，因此无法针对标签列和指标列进行计算，也就无法实现变量分箱计算

通过spark将样本**行列转置**之后，map过程每次传入的对象实际上就变成了样本的每列，因此可以计算变量分箱，再结合分布式任务，可以实现快速计算

```
257 @staticmethod
258 def online_data_trans_rdd(data, data_out, flag_name, not_var_list=[]):
259     """
260     """
261     Explanation
262
263     (1) 线上: 对spark里的dataframe进行行列转置, 变成rdd
264     (2) 线下: 或者对pandas里的dataframe进行转置
265     """
266     if type(data_out) == dataframe.DataFrame:
267         data_total = data.union(data_out)
268     else:
269         data_total = data
270     y = data_total.select(flag_name).toPandas()
271     n_sample = data.count()
272     data_total = data_total.select(list(set(data_total.columns) - set(not_var_list)))
273     columns = data_total.columns
274     # 获得变量列名, 用于标记spark每次map时计算的变量名称, 变量名称和变量在list中的序号一一对应
275     rdd_columns = sc.parallelize([tuple(columns)])
276     # 将变量名称和样本数据构造成rdd, 此时每次map的对象仍然是样本的一行
277     rdd_tuple = rdd_columns + data_total.rdd.map(tuple)
278     # 将n个样本n个指标拉平为m*n条数据, 每行数据为变量取值、变量名称、变量序号、并且通过zipwithindex从0到m*n-1对m*n条数据进行编号
279     rdd_add_index = rdd_tuple.zipWithIndex().flatMap(lambda x: [(x[1], k[0], k[1]) for k in enumerate(x[0])])
280     # 再根据变量序号进行分组, 将同一变量的所有取值合并到一起, 此时每次map的对象已经变为同一变量所有样本取值
281     rdd_groupby_index = rdd_add_index.map(lambda x: (x[1], (x[0], x[2]))) .groupByKey().sortByKey()
282     # 再根据样本index排序, 使得不同变量的取值一一对应
283     rdd_combine = rdd_groupby_index.map(lambda x: sorted(list(x[1]), key=lambda k:k[0]))
284     # 再通过map函数去掉index, 只保留变量名称和变量取值列表
285     rdd_transed = rdd_combine.map(lambda x: map(lambda y: y[1], x))
286     # 最后通过map函数将变量名称和变量取值列表分开
287     rdd_final = rdd_transed.map(lambda x: (x[0], x[1:]))
288     # 至此, 已经实现样本的行列转置, 得到我们需要的数据
289     return y, rdd_final, n_sample
290
291
292
293
294
295
296
297
```



原始样本只保留标签列和指标列



对变量分组的同时，根据分组结果将对应的变量进行WOE编码（主要应对于逻辑回归入模）

③ 样本woe转换

```
209 @staticmethod
210 def get_trans_woe_data(data_bin,data,factor_name):
211     """
212     Explanation
213
214     根据分bin结果将指标原始值转换为woe
215     1. 缺失值单独使用pandas.replace函数
216     2. 对于连续型指标, 使用pandas.cut函数
217     3. 对于离散型指标, 使用pandas.replace函数
218     """
219     data_select=data[[factor_name]]
220     tmp1=data_bin[data_bin['Bin_Type']=='Order_Bin']
221     tmp2 = data_bin[data_bin['Bin_Type'] != 'Order_Bin']
222     char_type=data_bin['Char_Type'][0]
223     if char_type=='numeric':
224         knot_list = [-numpy.inf] + list(map(lambda x: float(x.split(',')[1][-1]:-1]), tmp1['Bin']))
225         woe_list = list(tmp1['WOE'])
226         if len(tmp2)==0:
227             return pandas.DataFrame(pandas.cut(data_select[factor_name], bins=knot_list, labels=woe_list,include_lowest=True))
228         else:
229             single_dict=dict(zip(tmp2['Bin'],tmp2['WOE']))
230             single_bin_value, not_in_single = Base_Feature_Analyse.add_single_value(factor_name)
231             data1=data_select[(data_select[factor_name].isin(not_in_single)) | data_select[factor_name].isnull()]
232             data2=data_select.loc[data_select.index^data1.index]
233             data2[factor_name]=pandas.cut(data2[factor_name], bins=knot_list, labels=woe_list,include_lowest=True)
234             data2[factor_name]=data2[factor_name].astype(float)
235             for key in single_dict:
236                 data1=data1.replace(single_bin_value[key],single_dict[key])
237             if 'NA' in single_dict:
238                 data1=data1.fillna(single_dict['NA'])
239             trans_data=pandas.concat([data1,data2]).sort_index()
240             return trans_data
241     else:
242         single_dict = dict(zip(data_bin['WOE'],data_bin['Bin']))
243         if len(tmp2)==0:
244             for key in single_dict:
245                 data_select = data_select.replace(single_dict[key],key)
246             return data_select
247         else:
248             single_bin_value, not_in_single = Base_Feature_Analyse.add_single_value(factor_name)
249             for woe in single_dict:
```

```
250     if type(single_dict[woe])!=list:
251         if single_dict[woe] in single_bin_value:
252             if single_dict[woe]=='NA':
253                 data_select = data_select.fillna(woe)
254                 single_dict[woe]=single_bin_value[single_dict[woe]]
255     data_select = data_select.replace(single_dict[woe], woe)
256     return data_select
257
258
259
260
261
262
263
264
265
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
```

一、所有变量分组信息结果展示如下

Characterist	Description	Bin	WOE	IV_Bin	IV_Total	KS	#Obs	#Good	#Bad	%Obs	%Bad	%Cumulative_Bad	%Bad_Rate	%Cumulative_Bad_Rate
RH_out	变量含义	(-inf, 57.5]	0.83	0.05	0.06	0.17	844	772	72	9.70%	4.68%	4.68%	8.53%	8.53%
		(57.5, 63.0]	0.10	0.00	0.06	0.07	459	384	75	5.28%	4.88%	9.56%	16.34%	11.28%
		(63.0, 71.3333333333]	0.00	0.00	0.06	0.09	994	818	176	11.43%	11.44%	21.00%	17.71%	14.06%
		(71.3333333333, 88.5]	-0.03	0.00	0.06	0.07	3244	2657	587	37.30%	38.17%	59.17%	18.09%	16.42%
		(88.5, inf)	-0.15	0.01	0.06	0.05	3156	2528	628	36.29%	40.83%	100.00%	19.90%	17.68%
T2	变量含义	(-inf, 20.9725]	-0.12	0.01	0.13	0.13	5899	4748	1151	67.83%	74.84%	74.84%	19.51%	19.51%
		(20.9725, 22.4266666667]	-0.10	0.00	0.13	0.14	1419	1146	273	16.32%	17.75%	92.59%	19.24%	19.46%
		(22.4266666667, 23.29]	0.34	0.01	0.13	0.09	522	453	69	6.00%	4.49%	97.07%	13.22%	19.04%
		(23.29, inf)	1.35	0.11	0.13	0.18	857	812	45	9.85%	2.93%	100.00%	5.25%	17.68%
RH_3	变量含义	(-inf, 39.2285714286]	0.33	0.06	0.19	0.08	4976	4312	664	57.22%	43.17%	43.17%	13.34%	13.34%
		(39.2285714286, 42.1333333333]	-0.09	0.00	0.19	0.23	1707	1381	326	19.63%	21.20%	64.37%	19.10%	14.81%
		(42.1333333333, 44.2966666666]	-0.31	0.01	0.19	0.14	1102	853	249	12.67%	16.19%	80.56%	22.60%	15.92%
		(44.2966666667, 45.09]	-0.31	0.01	0.19	0.18	467	361	106	5.37%	6.89%	87.45%	22.70%	16.30%
		(45.09, inf)	-1.27	0.11	0.19	0.35	445	252	193	5.12%	12.55%	100.00%	43.37%	17.68%
T8	变量含义	(-inf, 21.1611111111]	-0.45	0.07	0.22	0.06	2659	1991	668	30.57%	43.43%	43.43%	25.12%	25.12%
		(21.1611111111, 22.6111111111]	-0.08	0.00	0.22	0.06	2615	2120	495	30.07%	32.18%	75.62%	18.93%	22.05%
		(22.6111111111, 24.0666666666]	0.32	0.02	0.22	0.04	2161	1869	292	24.85%	18.99%	94.60%	13.51%	19.57%
		(24.0666666667, 25.18]	0.95	0.06	0.22	0.14	825	762	63	9.49%	4.10%	98.70%	7.64%	18.38%
		(25.18, inf)	1.50	0.07	0.22	0.16	437	417	20	5.02%	1.30%	100.00%	4.58%	17.68%
Press_mm_hç	变量含义	(-inf, 752.15]	0.15	0.01	0.07	0.14	2726	2302	424	31.34%	27.57%	27.57%	15.55%	15.55%
		(752.15, 756.488888889]	0.14	0.00	0.07	0.09	1903	1603	300	21.88%	19.51%	47.07%	15.76%	15.64%
		(756.488888889, 764.133333333]	0.03	0.00	0.07	0.12	2888	2389	499	33.21%	32.44%	79.52%	17.28%	16.27%
		(764.133333333, 766.6]	-0.33	0.01	0.07	0.22	740	570	170	8.51%	11.05%	90.57%	22.97%	16.87%
		(766.6, inf)	-0.83	0.04	0.07	0.15	440	295	145	5.06%	9.43%	100.00%	32.95%	17.68%
RH_2	变量含义	(-inf, 33.5266666667]	0.51	0.01	0.03	0.23	445	394	51	5.12%	3.32%	3.32%	11.46%	11.46%
		(33.5266666667, 35.6266666666]	0.27	0.00	0.03	0.19	511	439	72	5.88%	4.68%	8.00%	14.09%	12.87%
		(35.6266666667, 37.6966666666]	0.09	0.00	0.03	0.11	1056	882	174	12.14%	11.31%	19.31%	16.48%	14.76%
		(37.6966666667, 41.718]	0.05	0.00	0.03	0.08	3380	2806	574	38.86%	37.32%	56.63%	16.98%	16.15%
		(41.718, inf)	-0.16	0.01	0.03	0.10	3305	2638	667	38.00%	43.37%	100.00%	20.18%	17.68%

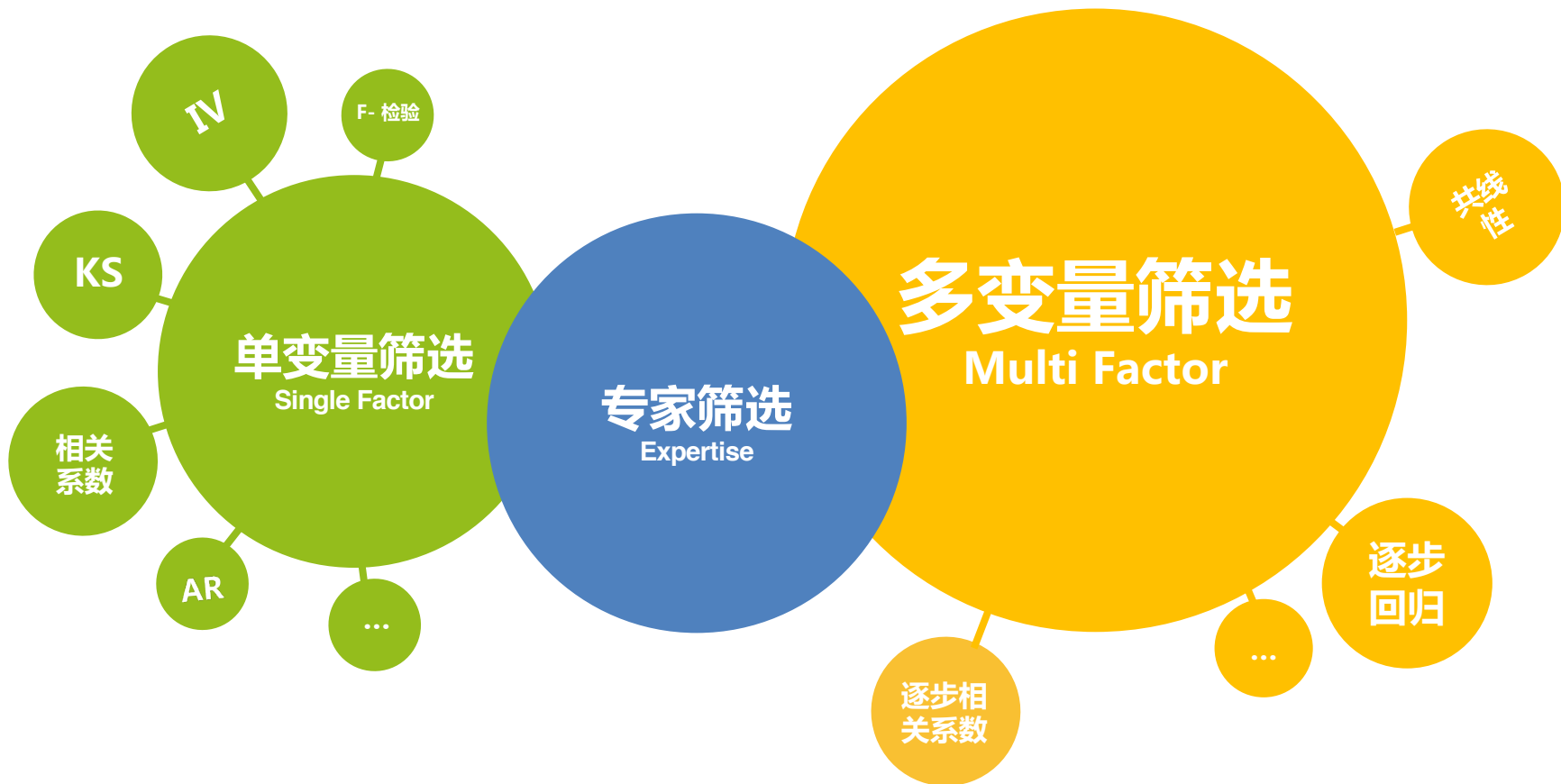
二、变量评价信息计算结果展示如下

14. 变量评价信息(Variables Effect)

Characteristic	Sample	AR	Corr_Pearson	Corr_Spearman	GINI	IV	KS	Nonempty_Bin_Num	Single_Threshold
samegps_xf_custnum_before1m_after1d_6geohash	Development Sample	0.28	0.22	0.17	0.28	0.29	0.21	4	43.22%
samegpscustnum1m3dec	Development Sample	0.26	0.20	0.15	0.26	0.28	0.20	5	3.03%
samegps_xf_over80_accountnum_before1m_after3d_6geohash	Development Sample	0.26	0.21	0.16	0.26	0.28	0.20	4	49.96%
samegpsapplnum1m3dec	Development Sample	0.26	0.20	0.14	0.26	0.28	0.20	5	2.95%
samegps_appl_count_1m	Development Sample	0.26	0.20	0.15	0.26	0.28	0.20	5	2.95%
samegpsrperiskpolicyrejcustnum1m3dec	Development Sample	0.26	0.20	0.14	0.26	0.26	0.19	5	4.30%
samegpsrejcustnum1m3dec	Development Sample	0.25	0.20	0.14	0.25	0.26	0.21	5	3.54%
samegps_xf_custnum_before1m_after3d_6geohash	Development Sample	0.27	0.21	0.16	0.27	0.25	0.20	4	34.92%
samegps_tx_or_xf_over80_accountnum_before1m_after3d_6geohash	Development Sample	0.28	0.21	0.16	0.28	0.25	0.23	5	20.85%
samegps_tx_or_xf_custnum_before1m_after3d_6geohash	Development Sample	0.27	0.21	0.15	0.27	0.25	0.22	5	20.72%
samegps_xf_over80_accountnum_before1m_after1d_6geohash	Development Sample	0.24	0.21	0.15	0.24	0.25	0.19	3	57.62%
samegpsapprcustnum1m3dec	Development Sample	0.26	0.20	0.14	0.26	0.25	0.20	5	11.92%
samegpsrejcustnum7d3dec	Development Sample	0.25	0.18	0.14	0.25	0.24	0.18	5	11.71%
samegpscustnum7d3dec	Development Sample	0.24	0.18	0.14	0.24	0.23	0.17	5	9.48%
samegps_tx_or_xf_over80_accountnum_before1m_after1d_6geohash	Development Sample	0.26	0.21	0.15	0.26	0.23	0.19	5	21.40%
samegpsapplnum7d3dec	Development Sample	0.24	0.18	0.14	0.24	0.23	0.17	5	9.27%
samegps_tx_or_xf_recordnum_before1m_after1d_6geohash	Development Sample	0.26	0.19	0.14	0.26	0.22	0.21	5	15.67%
samegpsrperiskpolicyrejcustnum7d3dec	Development Sample	0.25	0.18	0.14	0.25	0.22	0.19	5	12.64%
samegps_tx_or_xf_custnum_before1m_after1d_6geohash	Development Sample	0.26	0.21	0.15	0.26	0.22	0.20	5	21.82%
samegps_tx_or_xf_recordnum_before1m_after3d_6geohash	Development Sample	0.25	0.17	0.14	0.25	0.20	0.20	5	12.38%
samegps_xf_custnum_before1m_after1h_6geohash	Development Sample	0.20	0.18	0.14	0.20	0.19	0.18	3	66.39%
samegps_xf_recordnum_before1m_after1d_6geohash	Development Sample	0.24	0.16	0.14	0.24	0.18	0.19	5	43.22%
samegps_xf_recordnum_before1m_after3d_6geohash	Development Sample	0.23	0.13	0.13	0.23	0.17	0.17	5	34.92%
samegps_xf_custratio_before1m_after1d_6geohash	Development Sample	0.22	0.12	0.13	0.22	0.16	0.19	5	43.22%
samegps_tx_or_xf_recordnum_before1m_after1h_6geohash	Development Sample	0.21	0.17	0.12	0.21	0.16	0.16	5	22.45%
samegps_tx_custnum_before1m_after1d_6geohash	Development Sample	0.17	0.17	0.10	0.17	0.16	0.16	3	21.90%
samegps_first_xf_min_diff_before1m_after1d_6geohash	Development Sample	-0.17	-0.08	-0.10	-0.17	0.16	0.15	4	52.86%
samegps_tx_or_xf_custnum_before1m_after1h_6geohash	Development Sample	0.21	0.18	0.12	0.21	0.16	0.16	5	24.30%
samegps_tx_custnum_before1m_after3d_6geohash	Development Sample	0.17	0.17	0.10	0.17	0.15	0.17	5	22.87%
samegps_xf_over80_accountratio_before1m_after3d_6geohash	Development Sample	0.21	0.13	0.13	0.21	0.15	0.19	5	49.96%
samegps_xf_ratio_before1m_after1d_6geohash	Development Sample	0.09	0.08	0.06	0.09	0.15	0.12	3	52.86%
samegps_xf_ratio_before1m_after1h_6geohash	Development Sample	0.09	0.10	0.06	0.09	0.15	0.13	2	76.03%
samegps_xf_recordnum_before1m_after1h_6geohash	Development Sample	0.19	0.14	0.14	0.19	0.14	0.18	3	66.39%



变量筛选





多变量分析函数

- 在第二环节已经将单变量分析函数构造完成
- 此部分定义多变量分析函数
- **为何不将单变量分析函数在此部分定义**：所有变量一起分箱时一起计算单变量评价指标，节省时间



变量筛选流程

- **制定标准变量筛选流程**
- 将制定好的变量筛选流程函数化，并且针对每一种筛选方法配制一个控制参数，保证筛选流程可以灵活组合



自动筛选变量

- **单变量筛选和多变量筛选可以实现自动化**
- 专家筛选在很多场景不可缺少，但是专家筛选很难自动化
- 因此预留专家筛选缺口，保证单变量筛选和多变量筛选的自动化

① 多变量分析函数

```
19 class Base_Feature_Select_Method(object):
20
21     def __init__(self):
22         pass
23
24     @staticmethod
25     def eliminate_high_corr(train_data_woe, train_data_stat, not_var_list, corr_rate):
26         """
27         Explanation
28
29         筛选出相关系数小于指定阈值的变量
30         对于相关系数大于等于指定阈值的变量，只保留iv最大的变量
31
32         Parameters
33
34         train_data_woe:pandas.core.frame.DataFrame
35         | 训练集woe化样本
36
37         train_data_stat:pandas.core.frame.DataFrame
38         | 训练集指标表现
39
40         not_var_list:list
41         | 不参与筛选的指标列表
42
43         corr_rate:float
44         | 相关系数阈值
45
46         Return
47
48         corr_list:list
49         | 训练集相关系数小于指定阈值的指标名称List
50         """
51         var_list = list(set(train_data_stat['Characteristic'])-set(not_var_list))
52         '''计算所有指标之间的相关系数，并且转化为矩阵'''
53         train_corr = train_data_woe[var_list].corr().as_matrix()
54         n=len(var_list)
55         for i in range(n):
56             '''将矩阵的上半部分变为0，因为相关系数矩阵为对称矩阵，因此nxn矩阵只需要比较n*(n+1)/2即可，不需要n*n次'''
57             train_corr[i]=list(map(lambda x:0 if x>=i else train_corr[i][x],range(n)))
58         max_corr=1
59         while max_corr>=corr_rate:
```



```
60 print(max_corr)
61 n=train_corr.shape[0]
62 '''找到相关系数最大的行列索引'''
63 max_arg=numpy.argmax(train_corr)
64 max_row=max_arg//n
65 max_col=max_arg%n
66 iv_row=train_data_stat[train_data_stat['Characteristic']==var_list[max_row]]['IV'].max()
67 iv_col = train_data_stat[train_data_stat['Characteristic'] == var_list[max_col]]['IV'].max()
68 '''根据iv比较出需要删除的指标'''
69 del_max=max_row if iv_row<iv_col else max_col
70 train_corr=numpy.delete(train_corr,del_max,axis=1)
71 train_corr = numpy.delete(train_corr, del_max, axis=0)
72 var_list.remove(var_list[del_max])
73 max_corr=train_corr.max()
74 print(max_corr)
75 '''如此循环,直到所有指标的相关系数均小于既定的阈值'''
76 return var_list
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
00
86
87
88
89
90
91
```

优先级顺序	规则名称	规则详细逻辑	规则参数设置	规则状态	备注
CS01	单一阈值筛选	1、计算每个变量的每种取值占比 2、设定单一阈值阈值 3、遍历每个变量： 如果 $\max(\text{变量的不同取值占比}) > \text{单一阈值}$ ，则剔除变量；否则保留变量	单一阈值=0.9	启用	每一次筛选都是承接上一次筛选结果
CS02	最小变量分组数量	1. 变量分组默认分为5组，如果存在缺失值，则分6组，缺失值单独分为一组;分组要求，每组数量占比不得低于5% 2. 如果除去缺失值，剩下的无法分出大于等于“最小变量分组数量”的组数，则剔除变量;否则保留	最小变量分组数量=3	启用	
CS03	IV筛选	1、计算每个变量的最优分组 2、设定IV阈值 3、遍历每个变量： 如果 $\text{变量IV} < \text{单一阈值}$ ，则剔除变量；否则保留变量	IV阈值=0.02	启用	
CS04	PSI筛选	1. 如果建模样本的时间点跨度几个月，甚至更长，可以按月切分或者更长步长切分 2. 选择时间最早的一组作为基本样本，剩下的作为验证样本，计算组间的PSI，剔除PSI大于“PSI阈值”的变量	PSI阈值=0.05	启用	该部分视样本时间点实际情况而定
CS05	相关系数筛选	1、计算所有变量之间的相关性 2、设定相关系数阈值 3、遍历每个变量： 选择出与该变量相关系数大于相关系数阈值的所有变量，只保留这些变量中IV 最大额那个变量，其余变量删除	相关系数阈值=0.8	启用	
CS06	专家筛选	1. 专家结合实际业务场景，分析指标实际含义 2. 对于表现不错但是实际含义无法解释，或者分组无法解释的变量，剔除		启用	该部分视实际业务场景和指标含义而定
CS07	AR筛选	1、设定随机切分比例 2、随机切分数据集为训练集、测试集，计算训练集、测试集每个变量的AR值 3、遍历每个变量： 如果该变量在训练集、测试集上的 AR值负号相反，则剔除变量，否则保留	随机切分比例=6:4	未启用	
		1、计算所有变量之间的最大互信息数MIC			

CS08	MIC筛选	2、设定MIC阈值 3、遍历每个变量： 如果该变量的MIC<MIC阈值，则剔除；否则保留	MIC阈值=0.1	未启用	
CS09	皮尔逊卡方检验筛选	1、计算所有变量的皮尔逊卡方检验统计量的pvalue 2、设定置信水平alpha 3、遍历每个变量： 如果该变量的pvalue>alpha，则剔除；否则保留	alpha=0.05	未启用	
CS10	似然比检验筛选	1、计算所有变量的似然比检验统计量的pvalue 2、设定置信水平alpha 3、遍历每个变量： 如果该变量的pvalue>alpha，则剔除；否则保留	alpha=0.05	未启用	
CS11	F检验筛选	1、计算所有变量的F检验统计量的pvalue 2、设定置信水平alpha 3、遍历每个变量： 如果该变量的pvalue>alpha，则剔除；否则保留	alpha=0.05	未启用	
CS12	随机逻辑回归筛选	1、输入所有变量和标签 2、设定变量筛选惩罚系数、变量筛选阈值、随机建模种子 3、通过RLR随机逻辑回归建模，计算出通过筛选的变量	筛选惩罚系数=3 变量筛选阈值=0.1 随机建模种子=0	未启用	
CS13	逐步回归筛选	1、输入所有变量和标签 2、设定选入变量置信水平alpha1、剔除变量置信水平alpha2 3、通过逐步回归，每次对逻辑回归通过p检验的变量进行选入，对未通过p检验的变量进行剔除，逐步回归至结束	alpha1=0.05 alpha2=0.05	未启用	

2 变量筛选流程

```
65 class Base_Feature_Select_Rule(object):
66
67     @staticmethod
68     def initial_feature_select(data_original, data_stat, discrete_list=[], single_threshold=0.95, iv=0.02, corr=0.85, mic=None, min_bin_num=2):
69         """
70         Explanation
71         根据评价指标阈值，初步筛选出符合条件的指标
72         """
73         create_list=list(set(data_stat['Characteristic'])-set(discrete_list))
74         try:
75             discrete_data_stat=data_stat[(data_stat['Single_Threshold']!= 'nan') & (data_stat['Characteristic'].isin(discrete_list))]
76             create_data_stat = data_stat[(data_stat['Single_Threshold'] != 'nan') & (data_stat['Characteristic'].isin(create_list))]
77         except:
78             discrete_data_stat=data_stat[(data_stat['Single_Threshold'].notnull()) & (data_stat['Characteristic'].isin(discrete_list))]
79             create_data_stat = data_stat[(data_stat['Single_Threshold'].notnull()) & (data_stat['Characteristic'].isin(create_list))]
80         exclue_dic={'Single_Threshold':[]}
81         if single_threshold:
82             create_data_stat = create_data_stat[create_data_stat['Single_Threshold'] <= single_threshold]
83             discrete_data_stat = discrete_data_stat[discrete_data_stat['Single_Threshold'] <= single_threshold]
84             exclue_dic['Single_Threshold']=list(set(data_stat['Characteristic'])-set(create_data_stat['Characteristic'])-set(discrete_data_stat['Characteristic']))
85         if iv:
86             exclue_list=reduce(lambda x,y:x+y, exclue_dic.values())
87             create_data_stat = create_data_stat[create_data_stat['IV'] >= iv]
88             discrete_data_stat = discrete_data_stat[discrete_data_stat['IV'] >= iv]
89             exclue_dic['IV'] = list(set(data_stat['Characteristic']) - set(exclue_list) - set(
90                 create_data_stat['Characteristic']) - set(discrete_data_stat['Characteristic']))
91         if mic:
92             exclue_list = reduce(lambda x, y: x + y, exclue_dic.values())
93             create_data_stat = create_data_stat[create_data_stat['MIC'] >= mic]
94             discrete_data_stat = discrete_data_stat[discrete_data_stat['MIC'] >= mic]
95             exclue_dic['MIC'] = list(set(data_stat['Characteristic']) - set(exclue_list) - set(
96                 create_data_stat['Characteristic']) - set(discrete_data_stat['Characteristic']))
97         if min_bin_num:
98             exclue_list = reduce(lambda x, y: x + y, exclue_dic.values())
99             create_data_stat = create_data_stat[create_data_stat['Nonempty_Bin_Num']>=min_bin_num]
100             exclue_dic['Min_Bin_Num'] = list(set(data_stat['Characteristic']) - set(exclue_list) - set(
101                 create_data_stat['Characteristic']) - set(discrete_data_stat['Characteristic']))
102         if corr:
103             exclue_list = reduce(lambda x, y: x + y, exclue_dic.values())
104             corr_list=Base_Feature_Select_Method.eliminate_high_corr(data_original,create_data_stat, [], corr)
105             exclue_dic['Corr'] = list(set(data_stat['Characteristic']) - set(exclue_list) - set(corr_list) - set(discrete_data_stat['Characteristic']))
```

```

105     corr_list=Base_Feature_Select_Method.eliminate_high_corr(data_original,crete_data_stat, [], corr)
106     exclue_dic['Corr'] = list(set(data_stat['Characteristic']) - set(exclue_list) - set(corr_list) - set(discrete_data_stat['Characteristic']))
107     return corr_list+list(discrete_data_stat['Characteristic']),exclue_dic
108 final_list=list(crete_data_stat['Characteristic'])+list(discrete_data_stat['Characteristic'])
109 return final_list,exclue_dic
110
111 @staticmethod
112 def further_feature_select(flag_name,not_var_list,train_data_stat,train_data_woe,
113                           C=3,selection_threshold=0.1,random_state=0,p_value=0.05,rlr_rule=True,coef_rule=True,pvalue_rule=True):
114     """
115     Explanation
116     根据评价指标阈值，进一步筛选出符合条件的指标
117     """
118     final_list = list(set(train_data_stat['Characteristic'])-set(not_var_list))
119     if rlr_rule:
120         final_list = Base_Feature_Select_Method.eliminate_instabe_by_RLR(train_data_woe[final_list + [flag_name]],
121                                 not_var_list, C, selection_threshold,
122                                 random_state, flag_name)
123     if coef_rule:
124         final_list = Base_Feature_Select_Method.eliminate_positive_coefficient(train_data_woe[final_list + [flag_name]],
125                                     train_data_stat[
126                                         train_data_stat['Characteristic'].isin(
127                                             final_list)], flag_name)
128     if pvalue_rule:
129         final_list = Base_Feature_Select_Method.eliminate_high_pvalue_backword(train_data_woe[final_list + [flag_name]],
130                                     flag_name, p_value)
131     return final_list
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150

```